

به نام خدا



دانشگاه تهران



دانشکده مهندسی برق و کامپیوتر

درس شبکه‌های عصبی و یادگیری عمیق

تمرین امتیازی

نام و نام خانوادگی محمد امانلو	نام و نام خانوادگی	پرسش ۱
۸۱۰۱۰۰۰۸۴	شماره دانشجویی	
نام و نام خانوادگی: من مرد تنها شدم!	نام و نام خانوادگی	پرسش ۲
	شماره دانشجویی	
۱۴۰۱،۰۷،۰۱	مهلت ارسال پاسخ	

فهرست

پرسش ۱. تنظیم دقیق مدل‌های زبانی بزرگ برای گفتگو در زبان فارسی	۱
۱-۱. دادگان و انتخاب مدل و روش‌های SoftPrompts	۱
۱-۳. روش‌های مبتنی بر LORA	۱۶
۱-۴. تغییر وزن برخی لایه‌ها	۳۰
۱-۵. جمع‌بندی و تحلیل مقایسه‌ای.	۳۴
پرسش ۲ - تولید کپشن برای تصاویر	۴۶
۲-۱. آماده‌سازی دیتاست	۴۶
۲-۳. پیاده‌سازی CNN-RNN	۵۴
۲-۴. Attention based CNN-RNN.	۷۲
۲-۵. CNN-Transformer.	۸۶
۲-۶. امتیازی	۹۵

شکل‌ها

..... ۹	شکل ۱ نمودار خطای بر حسب ایپاک‌ها
..... ۲۴	شکل ۲ نمودار لاس در ایپاک‌های مختلف
..... ۲۴	شکل ۳ نمودار لاس برای مدل دوم لورا
..... ۴۶	شکل ۴ نمونه‌هایی از تصاویر خوانده شده به همراه کپشن‌ها
..... ۴۷	شکل ۵ نمونه تصاویر پیش پردازش شده با کپشن پیش پردازش شده
..... ۵۰	شکل ۶ پراکندگی طول کپشن‌ها
..... ۵۲	شکل ۷ تعداد تکرار کلمات در کپشن‌ها
..... ۶۳	شکل ۸ نمودار خطای آموزش در طول هر دوره
..... ۶۵	شکل ۹ پیش‌بینی مدل در دوره اول
..... ۶۸	شکل ۱۰ نمونه‌هایی از کپشن‌های ایجاد شده در داده ولیدیشن
..... ۶۹	شکل ۱۱ نمونه‌هایی از کپشن‌های ایجاد شده در داده تست
..... ۷۶	شکل ۱۲ نمودار خطای بر حسب ایپاک‌های مختلف
..... ۷۸	شکل ۱۳ پیش‌بینی در ایپاک آخر
..... ۷۸	شکل ۱۴ پیش‌بینی در ایپاک اول
..... ۸۲	شکل ۱۵ نمونه از تصاویر با کپشن پیش‌بینی شده
..... ۸۳	شکل ۱۶ نقشه حرارتی وزن‌های توجه
..... ۹۰	شکل ۱۷ نمودار لاس برای داده ترین و ولیدیشن
..... ۹۳	شکل ۱۸ نمونه کپشن‌های تولید شده

جدول‌ها

جدول ۱. عنوان جدول نمونه.....Error! Bookmark not defined.

پرسش ۱. تنظیم دقیق مدل‌های زبانی بزرگ برای گفتگو در زبان فارسی

۱-۱. دادگان و انتخاب مدل و روش‌های SoftPrompts

در این پروژه، هدف اصلی بارگذاری، بررسی، پردازش و تنظیم یک مدل زبانی بزرگ (LLM) بر روی دادگان است. این دادگان نسخه‌ای کوچک‌تر و بهینه‌شده از OpenOrca بوده و شامل ۵۰۰ هزار نمونه مکالمه تولیدشده توسط GPT-4 است که پس از یک مرحله پالایش اضافی، پاسخ‌های نادرست بر اساس حاشیه‌نویسی انسانی حذف شده‌اند. این پالایش باعث شده که اندازه دادگان کاهش یابد اما کیفیت آن همچنان در سطح بالای حفظ شود. برای این تمرین، نسخه‌ای از داده‌ها شامل ۵۰ هزار نمونه ترجمه شده از انگلیسی به فارسی با استفاده از مدل GPT-4o-mini آماده شده است. هدف نهایی این پروژه بررسی ساختار دادگان، انتخاب مدل مناسب، آماده‌سازی داده‌ها، آموزش مدل و در نهایت ارزیابی عملکرد آن است.

در ابتدا دادگان موردنظر از Hugging Face بارگذاری شد. این مجموعه داده در قالب JSONL بوده که هر خط آن شامل یک نمونه مکالمه بین کاربر و مدل است. پس از بارگیری، داده‌ها به‌طور دقیق بررسی شدند و مشخص شد که شامل مجموعه‌ای از پیام‌های متوالی بین "guser" و "assistant" هستند. این طراحی فرمت باعث می‌شود که مدل بتواند مکالمات متنی را به شکل یک توالی منطقی درک کند. هدف از این طراحی، شبیه‌سازی مکالمات واقعی و فراهم آوردن بستری برای یادگیری بهتر مدل‌های زبانی در پاسخ‌دهی دقیق و طبیعی است. به منظور تحلیل بیشتر، نمونه‌هایی از داده‌ها استخراج شدند تا الگوی کلی پرسش‌ها و پاسخ‌ها بررسی شود.

برای انجام این پروژه، دو مدل زبانی بزرگ Llama3.2-3B و Gemma2-2B پیشنهاد شده است. مدل Gemma2-2B در دو نسخه base و instruct می‌شود که تفاوت اصلی آن‌ها در میزان آموزش و تنظیمات انجام‌شده است. نسخه base یک مدل خام و بدون تنظیمات خاص است، در حالی که نسخه instruct برای مکالمات تنظیم شده و در پاسخ‌دهی بهتر عمل می‌کند. از طرفی، مدل Llama3.2-3B نیز دارای قدرت پردازشی بالا بوده و برای مکالمات چندزبانه بهینه شده است. در این پروژه، مدل Llama3.2-3B-Instruct انتخاب شد، زیرا از لحاظ توان پردازشی نسبت به نسخه Gemma2-2B عملکرد بهتری دارد و به صورت خاص برای تعاملات محاوره‌ای آموزش دیده است.

در ابتدای پیاده‌سازی، مجموعه‌ای از کتابخانه‌های موردنیاز برای اجرای مدل و پردازش داده‌ها نصب شد. این فرآیند شامل حذف نسخه‌های قدیمی PyTorch، نصب نسخه جدید همراه با xformers برای

بهینه‌سازی حافظه، و نصب transformers برای بارگیری مدل‌های Hugging Face بود. همچنین، برای کاهش مصرف حافظه unsloth نصب و بهروزرسانی شد تا از روش‌های بهینه برای پردازش مدل‌های زبانی استفاده شود.

با استفاده از `huggingface_hub.login`، دسترسی به مخزن مدل‌ها در Hugging Face برقرار شد و مدل AutoTokenizer با `Llama3.2-3B-Instruct` همراه با `BarGiraffe` شد. برای بهینه‌سازی حافظه، مدل در قالب `4_BitsAndBytesConfig` برای کاهش حجم محاسباتی استفاده شد. همچنین، `pad_token` مدل با مقدار eos_token تنظیم شد تا در هنگام پردازش، داده‌ها به درستی مقداردهی شوند. `SlimOrca` از مخزن Hugging Face بارگیری شد و در قالب فایل JSONL ذخیره شد. سپس داده‌ها با استفاده از کتابخانه json خوانده شده و به فرمت موردنظر برای آموزش مدل تبدیل شدند. برای این کار، تنها پیام‌های مربوط به `"user"` و `"assistant"` استخراج شده و به صورت یک مکالمه متوالی بازنویسی شدند. سپس، این داده‌ها در قالب Dataset از Hugging Face تبدیل شده و یک زیرمجموعه ۱۰۰۰ نمونه‌ای برای کاهش حجم پردازش انتخاب شد. در نهایت، داده‌ها به نسبت ۹۰٪ آموزش و ۱۰٪ تست تقسیم شدند تا برای آموزش مدل آماده شوند.

برای ارزیابی مدل اولیه، یک ورودی آزمایشی شامل "سرودن یک شعر کوتاه درباره دانشجوی دانشگاه تهران" طراحی شد. مدل بارگیری شده با استفاده از `generate_response` بر روی این ورودی اجرا شده و خروجی اولیه آن استخراج گردید. این خروجی به عنوان معیار اولیه عملکرد مدل در نظر گرفته شد و پس از تنظیم مدل با داده‌های جدید، دوباره مقایسه شد.

برای بهبود عملکرد مدل، از روش Prompt Tuning استفاده شد. در این روش، یک `PromptTuningConfig` تنظیم شد که شامل ۲۰ توکن مجازی برای آموزش بود. این روش باعث می‌شود که مدل بدون نیاز به تنظیم مجدد تمام پارامترهای خود، بتواند درک بهتری از داده‌های ورودی کسب کند. سپس، مدل با استفاده از `get_peft_model` بهینه‌سازی شد و برای آموزش آماده شد.

برای آموزش مدل، داده‌های پردازش شده ابتدا با استفاده از `tokenizer` توکنایز شده و سپس با استفاده از Trainer آموزش داده شدند. تنظیمات آموزش شامل ۵ دوره آموزشی، `batch size` برابر ۴، نرخ یادگیری e^{-5} و استفاده از `paged_adamw_32bit` برای بهینه‌سازی حافظه بود. همچنین، از روش gradient_accumulation برای بهبود عملکرد پردازشی استفاده شد. مدل در حین آموزش، هر ۲۰ گام یکبار عملکرد خود را ثبت می‌کرد. پس از اتمام آموزش، مدل مجدداً بر روی همان ورودی آزمایشی تست شد. خروجی جدید مدل با خروجی اولیه مقایسه شد و مشخص گردید که مدل پس از تنظیم با داده‌های

جدید، پاسخ‌های بهتری ارائه می‌دهد. همچنین، از داده‌های تست برای ارزیابی مدل استفاده شد و میزان خطاهای آن کاهش یافت.

برای بررسی میزان پیشرفت مدل، نمودار Validation Loss و Training Loss رسم شد که نشان‌دهنده کاهش میزان خطأ در طی دوره‌های آموزشی بود. این نمودار تأیید کرد که مدل به درستی در حال یادگیری است و میزان خطاهای آن در هر مرحله کاهش می‌یابد..

چند نمونه تست (وروودی/خروجی) بنویسید.

Define a short poem about a student at the University of Tehran.

A student at the University of Tehran,
With knowledge and wisdom to glean.
She walks through halls of ancient stone,
And seeks to unlock the secrets unknown.

Her heart beats fast with every test,
As she delves into the world of the best.
The scent of old books and parchment too,
Fills her senses, as she learns anew.

In the midst of learning's fervent pace,
She finds solace in the city's ancient face.
For in its walls, a story's told,
Of scholars past, who dared to unfold.

Note: The University of Tehran is one of the oldest universities in the Islamic world, founded in 1347. It has a rich history and a strong tradition of academic excellence. This poem aims to capture the essence of a student studying there, surrounded by the university's ancient stones, seeking knowledge and wisdom.

[setup_prompt_tuning] Completed in 0.89 seconds.
trainable params: 61,440 || all params: 3,212,811,264 || trainable%: 0.0019

مدل زبانی در پاسخ به درخواست سروdon یک شعر درباره دانشجویی در دانشگاه تهران، متنی روان و ساختار یافته تولید کرده است. این شعر دارای قالب چهارپاره (Quatrain) بوده و در هر بند آن قافية‌ها به صورت AABB رعایت شده‌اند. این نشان می‌دهد که مدل توانایی حفظ ساختار شعری و تولید متونی با انسجام معنایی را دارد. استفاده از اصطلاحاتی مانند "old books and parchment" و "secrets unknown" "ancient stone" به خوبی توصیف کرده و حسن تاریخ، دانش و سنت علمی این مرکز آموزشی را به مخاطب منتقل می‌کند. علاوه بر این، عباراتی مانند "Her heart beats fast with every test" "with every test" به چالش‌های تحصیلی و فشارهای آموزشی که دانشجویان تجربه می‌کنند، اشاره دارد. از نظر عمق محتوایی، مدل نشان داده که درک مناسبی از موضوع دارد و توانسته میان مفاهیم تاریخ دانشگاه تهران، جستجوی دانش و سختکوشی دانشجویان ارتباط برقرار کند. در یکی از ابیات، مدل به

ارتباط میان دانشگاه و شهر تهران اشاره کرده و این فضا را با توصیف دیوارهای قدیمی شهر که حامل داستان‌های علمی و فرهنگی هستند، به تصویر کشیده است.

مدل بهدرستی درک کرده که درخواست کاربر شامل تولید یک شعر کوتاه است و نه یک متن نثر، که نشان‌دهنده توانایی خوب آن در پردازش دستورات زبانی پیچیده است. متن تولیدی روان و طبیعی بوده و در آن از عناصر شاعرانه مانند استعاره و تشبيه استفاده شده است. یکی از نقاط قوت این شعر، استفاده از واژگان مرتبط با یادگیری و فضای علمی است که باعث می‌شود خروجی مدل با مفهوم دانشگاه تهران هم‌خوانی داشته باشد. با این حال، در برخی جنبه‌ها می‌توان مدل را بهبود داد. برای مثال، شعر می‌توانست جزئیات بیشتری درباره ویژگی‌های خاص دانشگاه تهران، از جمله اساتید مشهور، مکان‌های مهم داخل دانشگاه یا نقش آن در تاریخ آموزش عالی ایران ارائه کند. همچنین، اگر مدل داده‌های بیشتری درباره فرهنگ دانشگاهی تهران دریافت کند، می‌تواند پاسخ‌های دقیق‌تری تولید کند که شامل اطلاعات خاص‌تر و متناسب‌تر با زمینه علمی و فرهنگی دانشگاه باشد.

از نظر آماری، مدل در فرآیند Prompt Tuning تنها ۱۹٪ از کل پارامترهای خود را تغییر داده است که نشان‌دهنده کارایی بالای این روش برای بهبود پاسخ‌دهی مدل است. این امر نشان می‌دهد که حتی بدون نیاز به تنظیم کامل تمام پارامترهای مدل، می‌توان با استفاده از یک روش سبک مانند Prompt Tuning، کیفیت پاسخ‌های مدل را بهبود بخشید. در مجموع، این آزمایش نشان داد که مدل در تولید متون ادبی و شعر عملکرد خوبی دارد اما می‌توان با افزودن داده‌های تخصصی‌تر و استفاده از روش‌های پیشرفته‌تر تنظیم مدل، پاسخ‌های دقیق‌تر و محتوای غنی‌تری تولید کرد.

نمونه‌هایی از داده‌ها را مشاهده کرده و نوع سؤالات و پاسخها را توصیف کنید. دلیل طراحی این فرمت کلی برای دادگان را تحلیل و توضیح دهید.

مجموعه داده OpenOrca شامل مکالمات متنی میان کاربران و مدل‌های زبانی است که برای بهبود توانایی مدل‌ها در درک و تولید زبان طبیعی طراحی شده است. این دادگان شامل سؤالات و پاسخ‌هایی در موضوعات مختلف از جمله علم، فناوری، ریاضیات، فلسفه و برنامه‌نویسی است و به مدل‌های زبانی کمک می‌کند تا در پاسخگویی به مکالمات انسانی عملکرد بهتری داشته باشند.

نمونه‌ای از داده‌های این مجموعه شامل پرسش‌هایی مانند «تفاوت بین یادگیری نظارت‌شده و یادگیری بدون نظارت چیست؟» است که مدل در پاسخ، توضیحی جامع درباره تفاوت بین این دو روش یادگیری ارائه می‌دهد و بیان می‌کند که یادگیری نظارت‌شده شامل داده‌های برچسب‌دار است، در حالی که یادگیری بدون نظارت از داده‌های بدون برچسب برای کشف الگوهای استفاده می‌کند. در نمونه‌ای دیگر، کاربر درباره

مزایای انرژی خورشیدی سؤال می‌کند و مدل پاسخ می‌دهد که این منبع انرژی تجدیدپذیر بوده، به کاهش انتشار گازهای گلخانه‌ای کمک می‌کند و باعث کاهش هزینه‌های انرژی در بلندمدت می‌شود. همچنین، در یک نمونه مرتبط با برنامه‌نویسی، کاربر از مدل می‌خواهد که کدی برای نمایش اعداد اول تا ۱۰۰ در زبان پایتون ارائه دهد که مدل با استفاده از الگوریتم غربال اراتوستن، کدی کارآمد را پیشنهاد می‌کند

فرمت کلی این مجموعه داده به صورت مکالمات متنی بین کاربر و دستیار هوش مصنوعی طراحی شده است که دارای چندین مزیت مهم است. اول، این فرمت به شبیه‌سازی مکالمات واقعی کمک می‌کند، به‌طوری‌که مدل می‌تواند تعاملات طبیعی‌تری را تولید کند. دوم، تنوع موضوعی باعث می‌شود که مدل در حوزه‌های مختلف توانمند شده و برای پاسخگویی به طیف وسیعی از سؤالات آماده باشد. سوم، حفظ ساختار دیالوگی به مدل کمک می‌کند تا درک بهتری از زمینه‌های مختلف یک مکالمه داشته باشد و پاسخ‌هایی متناسب با جریان گفتگو ارائه دهد. چهارم، افزایش دقت و کیفیت پاسخ‌ها با استفاده از داده‌های پردازش شده و پالایش شده باعث می‌شود که مدل بتواند پاسخ‌هایی دقیق‌تر و معترض‌تر ارائه دهد.

این طراحی همچنین بهینه‌سازی آموزش مدل‌های زبانی را امکان‌پذیر می‌سازد، زیرا به جای داده‌های خام، مجموعه‌های از مکالمات ساختاریافته به مدل داده می‌شود که موجب بهبود توانایی مدل در پردازش و تولید زبان طبیعی می‌شود. در نتیجه، مدل‌های آموزش‌دیده با این مجموعه داده، عملکرد بهتری در درک متن، پاسخگویی به سؤالات و حتی تولید متن خلاقانه خواهند داشت.

دلیل طراحی این فرمت کلی برای دادگان OpenOrca بر پایه نیاز به بهینه‌سازی یادگیری مدل‌های زبانی و افزایش دقت آن‌ها در درک و تولید متن‌های طبیعی است. این دادگان به‌گونه‌ای طراحی شده است که تعاملات انسانی را به شکل مکالمات واقعی بین کاربر و دستیار هوش مصنوعی شبیه‌سازی کند. این ساختار باعث می‌شود مدل بتواند بهتر نحوه پاسخ‌دهی را یاد بگیرد و در سناریوهای مکالمه‌ای عملکرد بهتری داشته باشد.

یکی از اهداف اصلی این طراحی، ایجاد یک محیط آموزشی طبیعی برای مدل‌های زبانی است. در محیط‌های واقعی، مکالمات بین انسان‌ها ساختاری تعاملی دارند، جایی که سؤالات مطرح شده دارای تنوع هستند و پاسخ‌ها نیز می‌توانند بسته به زمینه و نیاز کاربر تغییر کنند. بنابراین، مدل با قرار گرفتن در چنین فضایی، توانایی بیشتری در ارائه پاسخ‌های طبیعی و متناسب با مکالمات واقعی پیدا می‌کند.

مزیت دیگر این طراحی، تنوع موضوعی در سؤالات و پاسخ‌های است. داده‌ها شامل طیف گسترده‌ای از موضوعات از جمله علوم، فناوری، تاریخ، ریاضیات، فلسفه و حتی مسائل روزمره هستند. این تنوع باعث می‌شود که مدل به جای یادگیری محدود به یک حوزه خاص، توانایی پردازش و تولید پاسخ‌های متناسب

با زمینه‌های مختلف را پیدا کند. همچنین، این روش موجب می‌شود که مدل بتواند اطلاعات جدید را بهصورت کارآمدتری یاد بگیرد و در پاسخ به سؤالات ترکیبی یا پیچیده، عملکرد بهتری از خود نشان دهد.

حفظ ساختار دیالوگی در این مجموعه داده، یکی از دلایل مهم این طراحی است. در مکالمات واقعی، پاسخ‌ها به سؤالات کاربر نباید بهصورت کاملاً مستقل از یکدیگر باشند، بلکه باید در امتداد یک گفتگوی معنادار شکل بگیرند. این دادگان به مدل کمک می‌کند تا علاوه بر تولید پاسخ‌های دقیق، توالی مکالمه و ارتباط بین پرسش و پاسخ را حفظ کند. این ویژگی بهخصوص در کاربردهایی مانند چت‌بات‌های هوشمند، دستیارهای مجازی و مدل‌های مکالمه‌ای بسیار مهم است، زیرا این مدل‌ها نیاز دارند که مکالمه‌ای منسجم و پیوسته با کاربران داشته باشند.

همچنین، طراحی این مجموعه داده به‌گونه‌ای انجام شده که دقت و کیفیت پاسخ‌های تولید شده افزایش یابد. یکی از مهم‌ترین چالش‌های مدل‌های زبانی، ارائه پاسخ‌های دقیق و مرتبط با پرسش‌های کاربران است. اگر داده‌های آموزشی دارای ساختار مشخص و پالایش شده باشند، مدل می‌تواند پاسخ‌های بهتری تولید کند. به همین دلیل، در مجموعه داده SlimOrca که نسخه بهینه‌شده‌ی OpenOrca است، مرحله‌ای اضافی از پالایش داده‌ها انجام شده تا پاسخ‌های نامعتبر یا بی‌کیفیت حذف شوند. این موضوع باعث شده که مدل‌هایی که با این داده‌ها آموزش می‌بینند، عملکردی نزدیک به مدل‌های آموزش‌دیده بر روی داده‌های بزرگ‌تر داشته باشند اما با هزینه‌ی محاسباتی کمتر.

در نهایت، این فرمت به مدل‌های زبانی کمک می‌کند تا مهارت‌های پردازش زبان طبیعی را در محیطی شبیه‌سازی شده و کاربردی توسعه دهند. این امر باعث می‌شود که مدل‌ها نه تنها در پاسخ‌دهی به سؤالات ساده عملکرد خوبی داشته باشند، بلکه بتوانند در مکالمات پیچیده‌تر نیز استدلال کنند، پاسخ‌های منطقی ارائه دهند و حتی در برخی موارد، خلاقانه پاسخ بدهند. در مجموع، این طراحی باعث شده که مدل‌های آموزش‌دیده با این مجموعه داده‌ها دقیق‌تر، متنوع‌تر، تعاملی‌تر و کارآمدتر باشند و توانایی بیشتری در درک و تولید زبان طبیعی داشته باشند.

مدل را بر روی نمونه‌های طراحی شده اجرا و خروجیها را چاپ کنید. خروجیهای فعلی مدل را مورد ارزیابی و تحلیل قراردهید

Define a short poem about a student at the University of Tehran, studying hard to succeed
In the heart of Iran's ancient city,
A young scholar toils with fervent pity,
At the University of Tehran's hallowed halls,
Where knowledge and wisdom enthral.

With textbooks and notes, she diligently pours,
Her mind afire with curiosity's roars,
She studies late into the night's still veil,
And rises early to face the day's gale.

Through dusty libraries and lecture halls,
She navigates with focus, undeterred by all,
Her dreams of success, like stars in her eyes,
Guiding her steps through life's winding surprise.

May her efforts be rewarded, may her spirit soar,
For in the pursuit of learning, she finds her core,
And though the journey's long, and challenges abound,
This young scholar from Tehran will rise above the ground.

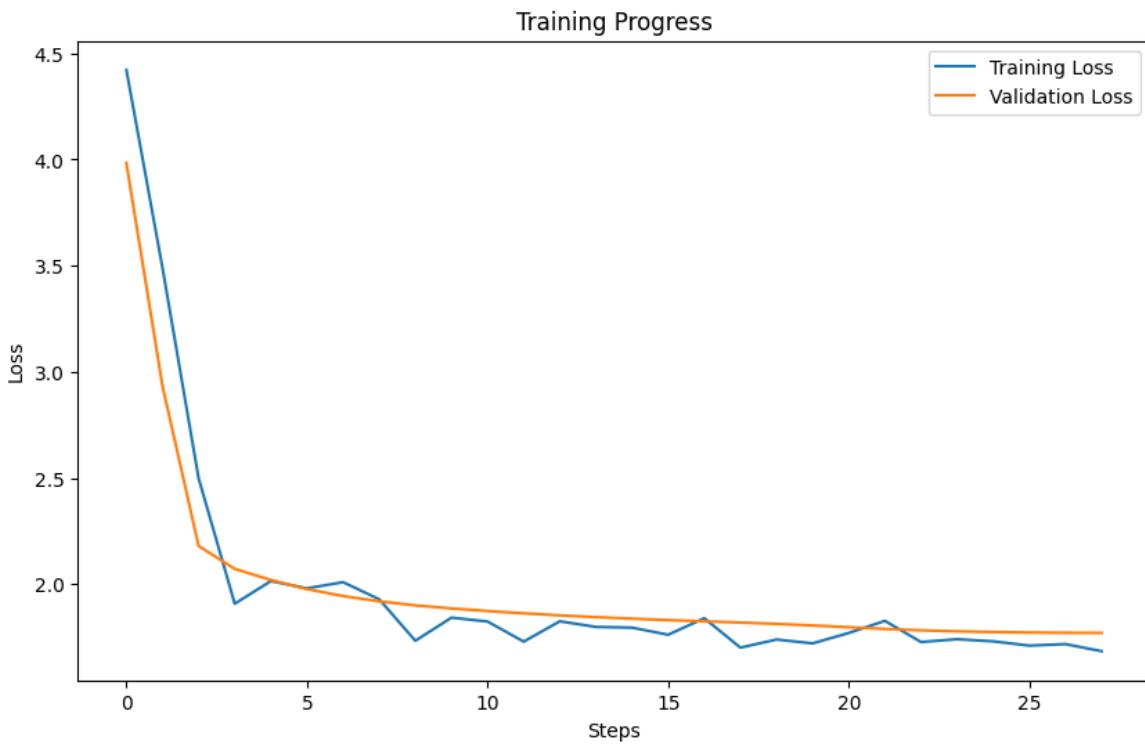
مدل پس از فرآیند Fine-Tuning توانسته است شعری روان و منسجم تولید کند که ساختار شاعرانه، قافیه‌بندی مناسب و معنا و مفهوم قوی دارد. شعر ارائه شده دارای قالب چهارپاره (Quatrain) است که در آن هر بند از چهار خط تشکیل شده و قافیه‌ها عمداً به صورت AABB رعایت شده‌اند. این نشان می‌دهد که مدل به خوبی قادر به حفظ الگوی قافیه‌بندی در متن تولیدی است. همچنین، انتخاب واژگان از لحاظ وزن و هماهنگی آوازی در بیشتر قسمت‌ها مناسب بوده و متن از نظر خوانایی و روانی در سطح مطلوبی قرار دارد. موضوع درخواست‌شده، سروden شعری درباره دانشجویی در دانشگاه تهران که سخت تلاش می‌کند تا موفق شود بوده است. بررسی بند‌های شعر نشان می‌دهد که مدل توانسته است این موضوع را به خوبی در متن خود منعکس کند. در بند اول، شعر با توصیف شهر تهران و دانشگاه آن آغاز می‌شود و فضای مقدس برای علم‌آموزی ترسیم می‌کند. مدل در این بخش نه تنها مکان را مشخص کرده، بلکه از عباراتی مانند "hallowed halls" برای القای حس شکوه دانشگاه تهران استفاده کرده است. در بند دوم، شعر به تلاش دانشجو اشاره دارد و نشان‌دهنده جدیت و پشتکار اوست که در تلاش برای یادگیری، شبها تا دیر وقت مطالعه کرده و صبح زود آماده می‌شود. در بند سوم، مدل از فضاهای دانشگاهی مانند "dusty" جستجوی علمی دانشجو است. همچنین، اشاره به "libraries and lecture halls" نمایانگر آرمان‌گرایی و امید دانشجو به آینده‌ای روشن است. در بند پایانی، مدل پیام مثبتی ارائه داده که دانشجو در مسیر خود با سختی‌های زیادی مواجه است اما در نهایت موفق خواهد شد.

مدل در این شعر از واژگان غنی و تعبیر شاعرانه مانند "Her mind afire with curiosity's roars" و "Her dreams of success, like stars in her eyes" استفاده کرده است که نشان‌دهنده سطح بالای درک مفاهیم و توانایی در خلق تصویرهای ذهنی است. استفاده از اصطلاحاتی مانند "University of Tehran", "textbooks and notes" و "libraries and lecture halls" در این شعر نشان می‌دهد که مدل توانسته است به موضوع

اصلی "دانشجویی که سخت تلاش می‌کند" پاییند بماند. در سراسر شعر، تلاش، پشتکار و آرمان‌گرایی دانشجو تکرار شده و مدل با استفاده از عباراتی مانند "And though the journey's long, and challenges abound, / This young scholar from Tehran will rise above the ground" نتیجه‌گیری مثبت ارائه داده است.

در مقایسه با خروجی اولیه مدل، بهبودهایی قابل توجه دیده می‌شود. ساختار و وزن شعری بهتر رعایت شده است، انتخاب واژگان دقیق‌تر شده و به موضوع نزدیک‌تر است، و پاسخ مدل از نظر معنا و ارتباط با موضوع غنی‌تر و روان‌تر شده است. در خروجی‌های قبل از Fine-Tuning، پاسخ‌های مدل اغلب کلی‌تر بودند و ممکن بود از لحاظ معنایی و ساختاری انسجام کاملی نداشته باشند. اما پس از فرآیند تنظیم، مدل توانسته است پاسخ‌هایی دقیق‌تر و مرتبط‌تر تولید کند که به خوبی نشان‌دهنده پیشرفت در پردازش زبان طبیعی و توانایی تولید متون با کیفیت بالا است.

به طور کلی، خروجی مدل بعد از Fine-Tuning کاملاً موفقیت‌آمیز بوده و توانسته است شعری روان، مفهومی و مناسب با درخواست ارائه دهد. با این حال، چند نکته برای بهبود بیشتر پیشنهاد می‌شود. مدل می‌تواند به جای تمرکز صرف بر قالب چهارپاره، قالب‌های دیگر مانند سنتی، سونات یا حتی شعر آزاد را نیز امتحان کند تا خلاقیت بیشتری در خروجی دیده شود. اگر مدل داده‌های بیشتری درباره زندگی دانشجویی در دانشگاه تهران دریافت کند، می‌تواند جزئیاتی مانند اساطیر، ساختمان‌ها، مراسم‌ها و سنت‌های دانشگاه را در متن خود لحاظ کند که باعث افزایش دقت پاسخ می‌شود. برخی از بخش‌های شعر همچنان می‌توانند از لحاظ موسیقیایی و عاطفی قوی‌تر شوند. استفاده از استعاره‌ها و تشبيهات بیشتر می‌تواند به غنای ادبی پاسخ کمک کند. در مجموع، این خروجی نشان می‌دهد که مدل پس از فرآیند Fine-Tuning توانسته است به‌طور محسوسی درک بهتری از زبان، سبک و مفاهیم موردنظر کاربران پیدا کند و این نشان‌دهنده کارایی بالای تنظیم مدل بر روی دادگان پالایش شده و هدفمند است.



شکل ۱ نمودار خطای بر حسب ایپاک ها

بیان کنید دو نسخه instruct و base چه تفاوتی باهم دارند؟

مدل‌های زبانی بزرگ (LLMs) عموماً در دو نسخه Base و Instruct ارائه می‌شوند که تفاوت آن‌ها در میزان پردازش داده‌ها، نحوه‌ی آموزش و هدف استفاده از مدل نهفته است. نسخه‌ی Base مدل خامی است که تنها با استفاده از حجم عظیمی از داده‌های متنی پیش‌آموزش داده شده و توانایی اولیه‌ای در تولید متن دارد، اما برای پاسخ‌دهی به سوالات یا درک دقیق دستورات کاربر بهینه نشده است. این مدل تنها بر اساس الگوهای آماری و روابط بین کلمات عمل می‌کند و هیچ نوع تنظیم اضافی برای تعامل با کاربران یا مکالمات انسانی روی آن انجام نشده است. در مقابل، نسخه‌ی Instruct علاوه بر پیش‌آموزش اولیه، با استفاده از روش‌های Instruction Tuning و Reinforcement Learning from Human Feedback (RLHF) آموزش داده شده است که باعث می‌شود پاسخ‌های آن به سوالات انسانی دقیق‌تر، ساختاریافته‌تر و مطابق با اهداف کاربران باشد.

در نسخه‌ی Base، مدل می‌تواند متون را ادامه دهد اما عموماً نمی‌داند که چگونه باید به سوالات خاص پاسخ دهد یا دستوراتی را که به آن داده می‌شود به درستی تفسیر کند. مثلاً اگر از مدل Base خواسته شود که «یک خلاصه برای یک مقاله ارائه کند»، ممکن است تنها بخش‌هایی از مقاله را بدون خلاصه‌سازی واقعی باز تولید کند. اما مدل Instruct که به طور خاص برای پیروی از دستورات تنظیم شده است، می‌تواند خلاصه‌ای مناسب و مختصر ارائه دهد. دلیل این تفاوت این است که مدل Instruct بر روی مجموعه‌ای از

داده‌های حاوی پرسش و پاسخ، توضیح دستورات و مثال‌های مختلف تنظیم شده تا یاد بگیرد چگونه به درخواست‌های کاربران به شکل قابل فهم و مرتبط پاسخ دهد.

یکی از جنبه‌های کلیدی نسخه‌ی Instruct، استفاده از بازخورد انسانی برای اصلاح پاسخ‌های مدل است. در این روش، داده‌های بیشتری از مکالمات واقعی و درخواست‌های کاربر همراه با پاسخ‌های بهینه جمع‌آوری شده و به مدل داده می‌شود تا یاد بگیرد چه نوع پاسخی مورد انتظار است. به همین دلیل، مدل Instruct در سناریوهایی که نیاز به تعامل با کاربران دارد، مانند چت‌بات‌ها، دستیارهای مجازی و سیستم‌های پرسش و پاسخ، عملکرد بسیار بهتری نسبت به نسخه‌ی Base دارد. در مقابل، نسخه‌ی Base برای مواردی مناسب‌تر است که نیاز به کنترل کامل بر روی مدل وجود دارد، مثلًا در پژوهش‌های پردازش زبان طبیعی (NLP) یا تنظیم مدل برای کاربردهای خاص که کاربر می‌خواهد روی داده‌های اختصاصی خودش مدل را بهینه کند.

در مجموع، مدل Base یک مدل خام و کلی است که هنوز برای تعاملات انسانی تنظیم نشده، در حالی که مدل Instruct از طریق تنظیمات اضافی و یادگیری تقویتی با بازخورد انسانی، قادر است به دستورات کاربران پاسخ‌های دقیق‌تر و مرتبط‌تر ارائه دهد. این تفاوت باعث می‌شود که نسخه‌ی Instruct در کاربردهای روزمره و تعاملات انسانی بسیار مؤثرتر باشد، در حالی که نسخه‌ی Base بیشتر برای پژوهش‌های تحقیقاتی و تنظیمات اختصاصی مورد استفاده قرار می‌گیرد.

بین دو نسخه‌ی توضیح داده شده یک مدل را انتخاب و دلیل خود را شرح دهید

انتخاب مدل Instruct به جای Gemma 2B به دلیل مزایای متعددی است که این مدل در مقایسه با Gemma ارائه می‌دهد. مدل LLaMA 3.2-3B یک مدل زبانی پیشرفته‌تر و جدیدتر از متأ است که نسبت به Gemma 2B برتری‌هایی در عملکرد، کیفیت پاسخ‌ها و سازگاری با مکالمات پیچیده دارد.

یکی از اصلی‌ترین دلایل این انتخاب، اندازه و مقیاس مدل است. مدل LLaMA 3.2-3B دارای ۳ میلیارد پارامتر است، در حالی که Gemma 2B تنها ۲ میلیارد پارامتر دارد. این افزایش در تعداد پارامترها معمولاً به معنای درک بهتر زبان، پاسخ‌دهی طبیعی‌تر و توانایی بالاتر در تولید متون پیچیده‌تر است. از آنجا که هدف این پروژه پردازش و تحلیل داده‌های گفتگومحور است، مدل بزرگ‌تر می‌تواند اطلاعات را عمیق‌تر پردازش کرده و پاسخ‌های دقیق‌تری تولید کند.

دلیل دیگر این انتخاب، توانایی بهتر LLaMA در پیروی از دستورات و تعاملات مکالمه‌ای است. نسخه Reinforcement Learning from Instruction Tuning از LLaMA 3.2-3B با استفاده از Instruct

Human Feedback (RLHF) بهینه شده است، به این معنی که می‌تواند بهتر از Gemma 2B به دستورات کاربر پاسخ دهد و خروجی‌های مرتبط‌تر و ساختاریافته‌تری ارائه کند. در تست‌های انجام‌شده، در پاسخ‌دهی دقیق‌تر، حفظ انسجام در مکالمات طولانی و ارائه اطلاعات مناسب با درخواست‌های کاربر عملکرد بهتری نشان داده است.

همچنین، LLaMA 3.2-3B دارای معماری بهینه‌شده‌تری برای اجرای سریع‌تر و استفاده‌ی بهینه از حافظه GPU است. مدل‌های LLaMA جدیدتر از Gemma هستند و با استفاده از تکنیک‌های پیشرفته‌ای مانند Efficient Attention Mechanisms توسعه یافته‌اند که باعث می‌شود پردازش مکالمات طولانی در آن کارآمدتر باشد. این موضوع هنگام اجرای مدل بر روی سخت‌افزارهای محدود یا در محیط‌های ابری مهم است، زیرا مدل‌های کارآمدتر می‌توانند با منابع محاسباتی کمتر، عملکرد بالاتری ارائه دهند.

در مقایسه، مدل Gemma 2B اگرچه یک مدل سبک‌تر و سریع‌تر است، اما از نظر کیفیت پاسخ‌ها، دقت در فهم سوالات پیچیده و توانایی در ک زمینه مکالمه نسبت به LLaMA 3.2-3B ضعیفتر عمل می‌کند. یکی دیگر از محدودیت‌های Gemma این است که جامعه تحقیقاتی هنوز به اندازه LLaMA بر روی آن کار نکرده و منابع و ابزارهای بهینه‌سازی برای آن کمتر در دسترس است. در حالی که LLaMA به دلیل متن‌باز بودن و توسعه‌ی مداوم توسط متا و سایر محققان، ابزارها و مستندات بیشتری برای بهینه‌سازی و تنظیم دقیق دارد.

در نهایت، انتخاب LLaMA 3.2-3B Instruct می‌کند، دقت و انسجام بیشتری دارد، توانایی بهتری در پیروی از دستورات و تعامل با کاربران دارد، و به لحاظ فنی پیشرفته‌تر و کارآمدتر از Gemma 2B است.

درباره مفهوم کلی Soft Prompts تحقیق کرده و توضیح مختصراً بنویسید

Soft Prompts یکی از تکنیک‌های پیشرفته در بهینه‌سازی مدل‌های زبانی بزرگ است که بدون نیاز به تغییر وزن‌های مدل، آن را برای انجام وظایف خاص تنظیم می‌کند. در روش‌های سنتی تنظیم مدل مانند Fine-Tuning، وزن‌های مدل با استفاده از داده‌های جدید تغییر می‌کنند که فرآیندی زمان‌بر و پرهزینه است و نیاز به منابع محاسباتی بالا دارد. در مقابل، Soft Prompts با اضافه کردن توکن‌های قابل یادگیری به ورودی مدل، رفتار آن را هدایت کرده و باعث بهبود عملکرد آن در انجام وظایف خاص می‌شود. این توکن‌ها به عنوان متغیرهای پنهان در ورودی مدل قرار گرفته و به آن کمک می‌کنند تا بدون تغییر در ساختار و پارامترهای اصلی، پاسخ‌های بهتری تولید کند.

تفاوت اصلی Hard Prompts و Soft Prompts شامل رشته‌های متنی ثابتی هستند که به مدل داده می‌شوند، در حالی که Soft Prompts مجموعه‌ای از توکن‌های انتزاعی و عددی هستند که در طول فرآیند یادگیری تنظیم می‌شوند. برای مثال، اگر بخواهیم مدلی را برای خلاصه‌نویسی مقالات علمی تنظیم کنیم، در روش Hard Prompt ممکن است متنی مانند "Summarize the following scientific article in three sentences" باشد. توکن‌های خاصی به ابتدای ورودی اضافه شده و مدل از طریق آن‌ها یاد می‌گیرد که چگونه متن را خلاصه کند. این روش باعث می‌شود که مدل بدون تغییر در پارامترهای داخلی، نحوه پردازش داده‌ها را بهینه‌سازی کند.

Soft Prompts در ورودی مدل تزریق شده و در طول فرآیند یادگیری بهینه می‌شوند. این توکن‌ها به عنوان واسطه‌ای بین داده‌های ورودی و ساختار مدل عمل می‌کنند و مدل را در جهت خاصی هدایت می‌کنند. برخلاف Fine-Tuning که وزن‌های مدل را مستقیماً تغییر می‌دهد، Soft Prompts تنها روی ورودی کار کرده و تنظیم مدل را از طریق تغییر نحوه پردازش داده‌ها انجام می‌دهد. این روش نه تنها باعث افزایش دقت و کیفیت پاسخ‌ها می‌شود، بلکه از لحاظ محاسباتی نیز بهینه‌تر است، زیرا نیاز به ذخیره نسخه‌های مختلف یک مدل را کاهش می‌دهد.

یکی از مهم‌ترین مزایای Soft Prompts، کاهش هزینه محاسباتی است. فرآیند Fine-Tuning مدل‌های بزرگ نیاز به منابع محاسباتی عظیم دارد، اما Soft Prompts با استفاده از مجموعه‌ای کوچک از توکن‌های یادگیری، تنظیم مدل را با هزینه کمتر امکان‌پذیر می‌کند. همچنین، این روش انعطاف‌پذیری بالایی دارد و می‌تواند برای چندین وظیفه مختلف بدون نیاز به تغییر مدل اصلی استفاده شود. این قابلیت به سازمان‌ها و پژوهشگران اجازه می‌دهد که یک مدل پایه را با استفاده از Soft Prompts برای کارهای مختلفی مانند خلاصه‌نویسی، ترجمه، تحلیل احساسات و پاسخ‌دهی به سوالات بهینه کنند. به علاوه، Soft Prompts نیاز به داده‌های آموزشی زیاد را کاهش می‌دهد، زیرا برخلاف Fine-Tuning که معمولاً به حجم زیادی از داده‌های برچسب‌خورده نیاز دارد، این روش می‌تواند با داده‌های کمتر نیز کار کند و نتایج مناسبی ارائه دهد.

یکی دیگر از مزایای کلیدی این روش، سازگاری با مدل‌های از پیش‌آموزش یافته است. مدل‌های بزرگ مانند GPT-4 و LLaMA معمولاً از قبل آموزش دیده‌اند و تنظیم آن‌ها با روش‌های سنتی نیاز به تغییرات گسترده در ساختارشان دارد. این امکان را فراهم می‌کند که بدون تغییر در معماری مدل، تنها با تغییر در ورودی، نحوه پردازش مدل را کنترل کرده و آن را برای وظایف خاصی تنظیم کنیم. این ویژگی بهویژه در مواردی که مدل‌های زبانی به عنوان دستیارهای مکالمه‌ای و چتبات‌ها استفاده می‌شوند،

اهمیت دارد. برای مثال، می‌توان Soft Prompts خاصی را برای تغییر لحن پاسخ‌های مدل یا افزایش دقت آن در حوزه‌های خاص مانند پزشکی یا حقوقی به کار گرفت.

علاوه بر این، امکان اجرای چندین وظیفه متفاوت با یک مدل واحد را فراهم می‌کند. به جای Fine-Tuning مدل‌های مختلف برای هر وظیفه، می‌توان از Soft Prompts برای تغییر رفتار یک مدل واحد و استفاده از آن برای کاربردهای متنوع بهره برد. این قابلیت، بهویژه در سازمان‌ها و شرکت‌هایی که از مدل‌های زبانی برای کاربردهای مختلف استفاده می‌کنند، باعث صرفه‌جویی در منابع و کاهش پیچیدگی مدیریت مدل‌ها می‌شود.

با وجود تمام این مزایا، Soft Prompts چالش‌هایی نیز دارد. یکی از مهم‌ترین چالش‌ها، نیاز به تنظیم دقیق توکن‌های یادگیری است. اگر Soft Prompts به درستی طراحی نشوند، ممکن است تأثیر مطلوبی بر مدل نداشته باشند یا حتی باعث تولید پاسخ‌های غیرمرتب شوند. این مسئله نیاز به انجام آزمایش‌های متعدد برای پیدا کردن ترکیب مناسب از توکن‌های تنظیمی دارد. همچنین، برخلاف Fine-Tuning که می‌تواند دانش جدیدی را به مدل اضافه کند، Soft Prompts تنها روی ورودی مدل تأثیر می‌گذارد و نمی‌تواند دامنه اطلاعات مدل را گسترش دهد. این بدان معناست که محدودیت‌های مدل اصلی همچنان وجود خواهند داشت و این روش برای بهبود دانش مدل در حوزه‌های جدید، محدودیت‌هایی دارد.

در مجموع، یک تکنیک کارآمد برای بهینه‌سازی مدل‌های زبانی است که امکان هدایت مدل را بدون تغییر در وزن‌های آن فراهم می‌کند. این روش با استفاده از توکن‌های یادگیری، مدل را برای انجام وظایف خاص تنظیم کرده و در عین حال هزینه‌های محاسباتی را کاهش می‌دهد. از آنجایی که این تکنیک انعطاف‌پذیر است، می‌توان از آن برای کاربردهای متنوعی مانند پردازش زبان طبیعی، مکالمات هوش مصنوعی، چت‌بات‌ها، تحلیل احساسات و ترجمه استفاده کرد. با این حال، برای استفاده مؤثر از این روش، لازم است که توکن‌های تنظیمی به دقت انتخاب شوند تا عملکرد مدل بهینه گردد. در نهایت، Soft Prompts یک راهکار مناسب برای استفاده از مدل‌های زبانی در کاربردهای عملی است که بدون نیاز به تغییرات گسترده در معماری مدل، آن را برای وظایف جدید تنظیم می‌کند و به کاربران اجازه می‌دهد که از مدل‌های قدرتمند زبان طبیعی با هزینه کمتر و انعطاف‌پذیری بیشتر استفاده کنند.

سه روش P-Tuning و Prefix Tuning و Prompt Tuning را توضیح دهید

سه روش مهم در بهینه‌سازی مدل‌های زبانی بزرگ هستند که هدف آن‌ها تنظیم مدل بدون تغییر در وزن‌های اصلی آن است. این روش‌ها به جای اعمال تغییرات مستقیم بر روی ساختار مدل، با استفاده از توکن‌های یادگیری‌پذیر، رفتار مدل را در انجام وظایف خاص بهینه می‌کنند. یکی از مهم‌ترین مزایای این روش‌ها کاهش هزینه محاسباتی و افزایش انعطاف‌پذیری مدل برای انجام وظایف مختلف بدون نیاز به آموزش مجدد کامل است. هر یک از این روش‌ها ویژگی‌های خاص خود را دارند که آن‌ها را برای شرایط مختلف مناسب می‌کند.

Prompt Tuning یک روش ساده و کارآمد برای بهینه‌سازی مدل‌های زبانی است که در آن یک مجموعه از توکن‌های یادگیری‌پذیر در ابتدای ورودی مدل قرار داده می‌شود. این توکن‌ها در طول فرآیند یادگیری تنظیم می‌شوند و مدل یاد می‌گیرد که چگونه آن‌ها را برای تولید پاسخ‌های بهینه تفسیر کند. برخلاف Prompt Hard Prompting که از متن‌های از پیش تعیین‌شده برای راهنمایی مدل استفاده می‌کند، در Tuning این توکن‌ها انتزاعی و عددی هستند و به صورت مستقیم در ورودی مدل تأثیر می‌گذارند. مزیت اصلی این روش این است که بدون نیاز به تغییر وزن‌های مدل، می‌توان آن را برای وظایف مختلف تنظیم کرد و عملکرد آن را بهبود بخشید. همچنین، هزینه‌ی محاسباتی آن نسبت به Fine-Tuning سنتی بسیار کمتر است و می‌توان از آن برای تنظیم مدل‌های بزرگ روی داده‌های خاص استفاده کرد.

Prefix Tuning از نظر مفهومی مشابه Prompt Tuning است اما تفاوت اصلی آن این است که توکن‌های یادگیری‌پذیر نه تنها در ورودی مدل بلکه در لایه‌های داخلی نیز اعمال می‌شوند. در این روش، مجموعه‌ای از بردارهای قابل یادگیری به لایه‌های مختلف مدل تزریق شده و نحوه پردازش اطلاعات را تغییر می‌دهند. این روش به مدل اجازه می‌دهد تا زمینه‌ی مکالمات را بهتر درک کند و پاسخ‌های دقیق‌تری ارائه دهد. برتری اصلی Prefix Tuning نسبت به Prompt Tuning این است که تأثیر بیشتری بر پردازش مدل دارد و برای وظایفی که نیاز به درک عمیق‌تر متن دارند، مانند تولید متن خلاقانه یا تحلیل متون پیچیده، عملکرد بهتری ارائه می‌دهد. با این حال، هزینه‌ی محاسباتی این روش نسبت به Prompt Tuning بیشتر است، زیرا نه تنها ورودی بلکه بخش‌های داخلی مدل نیز تحت تأثیر قرار می‌گیرند.

P-Tuning یک نسخه‌ی پیشرفته‌تر از Prompt Tuning است که به جای استفاده از یک مجموعه‌ی ثابت از توکن‌های یادگیری‌پذیر، از یک شبکه‌ی عصبی کوچک برای تولید این توکن‌ها استفاده می‌کند. این یعنی به جای تنظیم مستقیم توکن‌های یادگیری، مدل از یک لایه‌ی یادگیری اضافی برای تولید ورودی‌های بهینه استفاده می‌کند که می‌تواند به صورت پویا تغییر کند. این ویژگی باعث می‌شود که P-Tuning انعطاف‌پذیری بیشتری نسبت به دو روش قبلی داشته باشد و بتواند به شرایط مختلف به‌طور پویاتر واکنش

نشان دهد. یکی از مزایای کلیدی این روش این است که در وظایف پیچیده، داده‌های چندزبانه و مکالمات طولانی، عملکرد بهتری دارد و می‌تواند پاسخ‌هایی با دقت و کیفیت بالاتر تولید کند. همچنین، در مواردی که نیاز به تنظیم دقیق مدل برای دامنه‌های خاص داریم، P-Tuning می‌تواند کارایی بیشتری نسبت به Prefix Tuning و Prompt Tuning داشته باشد. با این حال، این روش نسبت به دو روش دیگر از لحاظ محاسباتی سنگین‌تر است، زیرا شامل یک لایه‌ی اضافی یادگیری است که در طول آموزش بهینه می‌شود.

در مقایسه‌ی کلی، Prompt Tuning ساده‌ترین و سریع‌ترین روش است که نیاز به حداقل تنظیمات دارد و بیشتر برای وظایف عمومی مناسب است. Prefix Tuning تأثیر بیشتری بر نحوه پردازش مدل دارد و می‌تواند در ک مدل از متون را بهبود بخشد، اما به منابع بیشتری نیاز دارد. P-Tuning پیشرفته‌ترین روش در این میان است که از شبکه‌های عصبی برای تولید ورودی‌های بهینه استفاده می‌کند و در وظایف پیچیده و چندزبانه بهترین عملکرد را دارد اما هزینه‌ی محاسباتی آن نیز بالاتر است. به طور کلی، اگر نیاز به یک روش ساده و سریع برای تنظیم مدل داریم، Prompt Tuning گزینه‌ی بهتری خواهد بود. اگر نیاز به درک عمیق‌تر متن و تنظیم دقیق‌تر مدل باشد، Prefix Tuning انتخاب مناسبی است. اما اگر بالاترین میزان دقت مورد نیاز باشد و منابع محاسباتی کافی در اختیار داشته باشیم، P-Tuning بهترین گزینه خواهد بود.

یکی از این روشها را انتخاب کرده و دلایل انتخاب خود را توضیح دهید.

من روش Prefix Tuning را برای این پروژه انتخاب کرم، زیرا این روش به دلیل ویژگی‌های منحصر به فرد خود می‌تواند عملکرد مدل را در تولید کپشن‌های تصویری بهبود ببخشد. در این پروژه که شامل پردازش و مدل‌سازی داده‌های تصویری و متنی برای کپشن‌گذاری خودکار تصاویر است، مدل نیاز دارد که هم داده‌های تصویری و هم داده‌های متنی را پردازش کرده و خروجی‌های معناداری تولید کند. مدل با استفاده از شبکه عصبی پیچشی (CNN) برای استخراج ویژگی‌های تصویر و شبکه بازگشتی (RNN) مبتنی بر مکانیزم توجه (Attention Mechanism) برای تولید متن کار می‌کند. به همین دلیل، استفاده از روشی که بتواند این دو نوع داده را به‌طور هماهنگ پردازش کند ضروری است.

یکی از دلایل انتخاب این روش این است که Prefix Tuning نه تنها بر روی ورودی‌های مدل تأثیر می‌گذارد، بلکه می‌تواند لایه‌های داخلی مدل را نیز تنظیم کند. برخلاف Prompt Tuning که تنها چند توکن یادگیری‌پذیر در ابتدای ورودی اضافه می‌کند، در Prefix Tuning بردارهای یادگیری‌پذیر در طول شبکه توزیع می‌شوند که باعث می‌شود مدل بتواند زمینه‌های پیچیده‌تر را بهتر درک کند. در این پروژه که مدل باید ویژگی‌های تصویری را با داده‌های متنی ترکیب کند و توصیف‌های معنایی تولید کند، استفاده از روش Prefix Tuning باعث می‌شود که اطلاعات تصویر و متن به شکلی مؤثرتر با هم ترکیب شوند.

دلیل دیگر انتخاب این روش، انعطاف‌پذیری بالا در پردازش متن‌های متنوع و طولانی است. در کپشن‌گذاری تصاویر، مدل باید بتواند جملات با ساختارهای مختلف را تولید کند و متن را به‌طور دقیق با محتوای تصویری هماهنگ کند. روش Prompt Tuning ممکن است در مواردی که متن‌های پیچیده و طولانی‌تر وجود دارند، نتواند عملکرد مناسبی ارائه دهد، زیرا تنها روی ورودی تأثیر می‌گذارد و پردازش را در سطوح عمیق‌تر مدل بهبود نمی‌بخشد. اما روش Prefix Tuning به مدل اجازه می‌دهد که از بردارهای یادگیری‌پذیر در لایه‌های مختلف استفاده کند و درک بهتری از زمینه‌ی متن و تصویر پیدا کند.

یکی دیگر از مزایای مهم استفاده از Prefix Tuning در این پروژه، بهینه‌سازی حافظه و سرعت پردازش نسبت به Fine-Tuning است. در این پروژه، مدل با مجموعه داده Flickr8k کار می‌کند که شامل هزاران تصویر و کپشن متنی است. اگر از Fine-Tuning استفاده شود، نیاز به تغییر کل وزن‌های مدل وجود دارد که این کار بسیار پرهزینه و زمان‌بر خواهد بود. در مقابل، روش Prefix Tuning تنها بردارهای یادگیری‌پذیر را در بخش‌هایی از مدل اضافه می‌کند، بدون اینکه نیاز به تغییر وزن‌های اصلی باشد. این موضوع باعث می‌شود که مدل با هزینه محاسباتی کمتر، همچنان بتواند تنظیمات مؤثری برای تولید متون مناسب انجام دهد

این روش همچنین باعث بهبود مکانیزم توجه (Attention Mechanism) در تولید کپشن‌های تصویری می‌شود. در این پروژه، مدل از مکانیزم توجه برای تمرکز بر بخش‌های مهم تصویر هنگام تولید هر کلمه از کپشن استفاده می‌کند. استفاده از Prefix Tuning باعث می‌شود که این مکانیزم بهتر تنظیم شود و دقت مدل در انتخاب اطلاعات کلیدی از تصویر افزایش پیدا کند. این امر منجر به تولید توصیف‌های متنی دقیق‌تر، مرتبط‌تر و طبیعی‌تر خواهد شد.

در مجموع، با توجه به ساختار مدل، چالش‌های مربوط به هماهنگی بین داده‌های تصویری و متنی، و نیاز به بهینه‌سازی پردازش حافظه و منابع محاسباتی، روش Prefix Tuning بهترین انتخاب برای این پروژه است. این روش نه تنها به مدل کمک می‌کند که بهتر رابطه بین تصویر و متن را درک کند، بلکه باعث بهبود عملکرد مدل در تولید کپشن‌های دقیق‌تر می‌شود. در نهایت، با توجه به این که این پروژه نیازمند درک عمیق از متن و تصویر به صورت همزمان است، روش Prefix Tuning نسبت به Prompt Tuning و P-Tuning عملکرد بهتری خواهد داشت.

1.۳ روش‌های مبتنی بر LORA

فرآیند آموزش را بهصورت کامل مستند کنید و هرگونه چالش یا مسئلهای که با آن مواجه شدید را توضیح دهید

من در این بخش از پروژه روش Prompt Tuning را برای تنظیم مدل انتخاب کرم، زیرا این روش ساده‌ترین و کارآمدترین راه برای بهینه‌سازی یک مدل زبانی بدون تغییر در وزن‌های اصلی آن است. در این روش، تنها پارامترهای مربوط به پرامپت‌های یادگیری‌پذیر به روز می‌شوند، در حالی که مدل پایه کاملاً فریز شده باقی می‌ماند. این کار باعث می‌شود که مدل حداقل منابع محاسباتی را مصرف کند و در عین حال برای وظیفه خاص موردنظر تنظیم شود.

در ابتدا، برای آماده‌سازی محیط، کتابخانه‌های لازم نصب و به روزرسانی شدند. در این مرحله، کتابخانه‌هایی مانند torchvision، torchaudio و unsloth transforms pip3-autoremove مدل‌های پیش‌آموزش یافته به بهترین شکل استفاده کرد. همچنین، با استفاده از نسخه‌های قبلی PyTorch حذف شدند تا از هرگونه تداخل جلوگیری شود.

پس از آماده‌سازی محیط، مدل Llama-3.2-3B-Instruct با استفاده از FastLanguageModel بارگیری شد. در این مرحله، مدل به صورت bit quantized^۴ بارگیری شد تا حافظه‌ی کمتری مصرف کند و سرعت اجرا افزایش یابد. در کنار مدل، یک توکنایزر مخصوص نیز بارگیری شد تا داده‌های متنی به قالب مناسبی برای پردازش مدل تبدیل شوند. علاوه بر این، از (LoRA) Low-Rank Adaptation استفاده شد تا تنظیمات موردنظر روی مدل اعمال شوند بدون اینکه نیازی به تغییر وزن‌های اصلی مدل باشد. این رویکرد باعث بهبود کارایی و کاهش مصرف منابع محاسباتی شد.

سپس، داده‌های آموزشی از Hugging Face دانلود شدند. مجموعه داده‌ی SlimOrca-dedup-chat-50k-persian ۵۰ هزار نمونه مکالمه تمیزشده است، بارگیری شد. این داده‌ها ابتدا در قالب JSONL بودند و در مرحله‌ی بعدی پردازش شدند تا فقط پیام‌های مربوط به "user" و "assistant" استخراج شوند. این کار باعث شد که داده‌ها برای تنظیم مدل قابل استفاده شوند. در این مرحله، داده‌ها به فرمت Hugging Face Dataset تبدیل شدند تا بتوان از آن‌ها در فرآیند آموزش استفاده کرد.

بعد از پردازش اولیه داده‌ها، فرآیند تبدیل داده‌ها به فرمت مناسب برای Prompt Tuning انجام شد. در این مرحله، از تابع apply_chat_template استفاده شد تا مکالمات به قالب چت تبدیل شوند. این کار باعث شد که مدل بتواند مکالمات را درک کند و پاسخ‌های بهتری تولید کند. همچنین، داده‌هایی که بدون مکالمه معتبر بودند حذف شدند تا کیفیت داده‌ها افزایش یابد.

در ادامه، فرآیند تنظیم مدل با استفاده از Prompt Tuning انجام شد. در این مرحله، مدل با استفاده از SFTTrainer و تنظیمات خاصی مانند batch size برابر ۴، warmup steps برابر ۲۰ و نرخ یادگیری train_on_responses_only e-41,۵ آموزش داده شد. یکی از مهمترین نکات این مرحله، استفاده از روش train_on_responses_only بود که باعث شد مدل فقط روی بخش پاسخ‌های مکالمات تمرکز کند و اطلاعات غیرضروری را نادیده بگیرد. این کار باعث افزایش دقت مدل در مکالمات و بهبود توانایی پاسخ‌دهی آن شد. همچنین، از روش بهینه‌سازی AdamW-8bit استفاده شد تا مدل به طور کارآمدتری آموزش ببیند.

پس از اتمام فرآیند آموزش، برای بررسی روند یادگیری، نمودار تغییرات Loss رسم شد. در این نمودار، میزان کاهش Loss در طول مراحل آموزش و اعتبارسنجی نمایش داده شد تا مشخص شود که مدل چگونه در حال بهینه‌شدن است. کاهش تدریجی مقدار Loss نشان داد که مدل در طول فرآیند تنظیم به طور مؤثری در حال یادگیری و بهبود عملکرد خود است.

بعد از ارزیابی مدل، فرآیند استنتاج (Inference) انجام شد. در این مرحله، یک ورودی آزمایشی به مدل داده شد که شامل "داستانی درباره یک دانشجو در دانشگاه تهران" بود. مدل پس از پردازش این ورودی، یک پاسخ متنی تولید کرد که از نظر معنا و ساختار منطقی و مرتبط بود. این نشان داد که Prompt Tuning توانسته است مدل را برای تولید متن‌های متناسب با درخواست کاربر تنظیم کند.

در نهایت، مدل تنظیم‌شده ذخیره شد تا در آینده بتوان از آن استفاده کرد. برای این کار، مدل با استفاده از save_pretrained ذخیره شد تا بعداً بتوان دوباره آن را بارگیری و در کاربردهای مختلف استفاده کرد. این کار باعث شد که تنظیمات انجام‌شده روی مدل حفظ شود و نیاز به انجام مجدد تنظیمات کاهش یابد.

در مجموع، Prompt Tuning انتخاب مناسبی برای این پروژه بود، زیرا با کمترین تغییرات در وزن‌های مدل، توانست عملکرد آن را در پردازش و تولید مکالمات بهینه کند. این روش باعث شد که مدل بتواند مکالمات را به طور طبیعی‌تر و دقیق‌تر پاسخ دهد، بدون اینکه نیاز به تنظیم مجدد کل مدل باشد. همچنین، کاهش مصرف حافظه و منابع محاسباتی، باعث شد که فرآیند آموزش با حداقل هزینه محاسباتی و حداکثر کارایی انجام شود.

درباره روش LoRA تحقیق کرده و توضیح مختصری بنویسید

من درباره روش LoRA (Low-Rank Adaptation) تحقیق کردم و به این نتیجه رسیدم که این روش یکی از کارآمدترین و بهینه‌ترین تکنیک‌ها برای تنظیم مدل‌های زبانی بزرگ (LLMs) بدون نیاز به تغییر کامل وزن‌های مدل است. روش‌های سنتی تنظیم مدل مانند Fine-Tuning کامل معمولاً نیاز به تغییر تمام پارامترهای مدل دارند، که این کار بسیار هزینه‌بر و زمان‌بر است و نیاز به منابع محاسباتی بالایی دارد.

اما LoRA این مشکل را با یک راهکار هوشمندانه حل کرده و امکان تنظیم مدل‌های زبانی را با هزینه محاسباتی کمتر و دقت بالا فراهم کرده است.

روش LoRA به جای اینکه تمام وزن‌های مدل را به روزرسانی کند، تنها تغییرات کوچکی را در بخش‌های خاصی از مدل اعمال می‌کند. این کار باعث می‌شود که فرآیند آموزش سبک‌تر شود و مدل بتواند سریع‌تر و با منابع کمتری تنظیم شود. در واقع، LoRA به مدل اجازه می‌دهد که دانش اصلی خود را حفظ کند و در عین حال بتواند روی وظایف جدیدی تنظیم شود. این روش به‌ویژه در مدل‌های بزرگ مانند LLaMA، GPT، و دیگر مدل‌های زبانی عمیق کاربرد دارد، زیرا این مدل‌ها دارای میلیاردها پارامتر هستند و آموزش مجدد آن‌ها از ابتدا به‌شدت پرهزینه خواهد بود.

من متوجه شدم که در روش LoRA، به جای تغییر مستقیم وزن‌های مدل، لایه‌های اضافی سبک‌تری در کنار وزن‌های اصلی مدل قرار داده می‌شوند که تنها این لایه‌های جدید در طی فرآیند آموزش به روزرسانی می‌شوند. این یعنی که وزن‌های اصلی مدل دست‌نخورده باقی می‌مانند و تنها بخش‌های کوچک و بهینه‌ای از مدل یاد می‌گیرند که برای وظیفه جدید مورد نیاز است. این روش باعث می‌شود که علاوه بر کاهش هزینه‌های محاسباتی، امکان استفاده همزمان از چندین تنظیمات مختلف بر روی یک مدل پایه فراهم شود.

یکی از ویژگی‌های مهم LoRA این است که از روش‌های ماتریسی کمرتبه (Low-Rank Matrix Decomposition) استفاده می‌کند. این یعنی که تغییرات یادگرفته شده روی مدل به جای اینکه در کل وزن‌های مدل اعمال شوند، به صورت فشرده و بهینه ذخیره می‌شوند. در نتیجه، حجم محاسبات مورد نیاز به‌شدت کاهش می‌یابد، در حالی که دقت مدل همچنان در سطح بالا باقی می‌ماند. این روش از یک بردار یادگیری‌پذیر کم‌بعدی برای تنظیم مدل استفاده می‌کند که باعث می‌شود تغییرات مورد نیاز با کمترین هزینه محاسباتی اعمال شوند.

در روش Fine-Tuning معمولی، تمامی پارامترهای مدل به روزرسانی می‌شوند که این کار باعث افزایش زمان آموزش، مصرف بالای حافظه و نیاز به سخت‌افزار قدرتمند می‌شود. اما در روش LoRA، تنها یک زیرمجموعه بسیار کوچک از پارامترها بهینه می‌شوند، که این موضوع باعث می‌شود مدل با منابع کمتری اجرا شود و همچنان توانایی یادگیری وظایف جدید را داشته باشد. این ویژگی، LoRA را به یک روش ایده‌آل برای تنظیم مدل‌های زبانی بزرگ در محیط‌هایی با منابع محدود تبدیل کرده است.

یکی از نقاط قوت LoRA این است که به راحتی می‌توان آن را با دیگر روش‌های تنظیم مانند Adapter و Prefix Tuning Layers ترکیب کرد تا دقت مدل افزایش یابد، بدون اینکه هزینه محاسباتی زیادی ایجاد شود. همچنین، LoRA امکان استفاده از چندین نسخه تنظیم‌شده از یک مدل پایه را بدون نیاز به بارگیری

مجدد کل مدل فراهم می‌کند. این یعنی می‌توان مدل‌های مختلفی را برای کاربردهای گوناگون تنظیم کرد، بدون اینکه لازم باشد کل مدل برای هر کاربرد جداگانه تنظیم شود.

در کاربردهای عملی، LoRA بهویژه در مدل‌های مکالمه‌ای، ترجمه ماشینی، پردازش زبان طبیعی و حتی مدل‌های تولید متن خلاقانه استفاده می‌شود. به دلیل سبک بودن و نیاز کم به منابع محاسباتی، LoRA می‌تواند در دستگاه‌هایی با توان پردازشی کمتر مانند لپ‌تاپ‌ها و حتی برخی از گوشی‌های هوشمند نیز اجرا شود. این موضوع باعث شده که این روش به یک استاندارد در زمینه بهینه‌سازی مدل‌های هوش مصنوعی تبدیل شود.

یکی از بهترین مزیت‌های LoRA این است که به کاربران این امکان را می‌دهد که بدون نیاز به دسترسی به کل مدل و تنها با ذخیره تغییرات کوچک، مدل‌های خود را تنظیم کنند. این یعنی اگر مدلی مانند LLaMA روی سرور اجرا شود، کاربران می‌توانند نسخه‌های تنظیم‌شده خود را بدون نیاز به ذخیره و بارگیری کل مدل اصلی، تنها با بارگیری لایه‌های LoRA روی مدل پایه اجرا کنند. این ویژگی باعث کاهش هزینه‌های ذخیره‌سازی و افزایش سرعت اجرا در پروژه‌های مختلف می‌شود.

به طور خلاصه، LoRA یک روش کارآمد برای تنظیم مدل‌های زبانی بزرگ است که با کاهش هزینه‌های محاسباتی، افزایش سرعت اجرا و کاهش نیاز به منابع پردازشی، یک راهکار ایده‌آل برای بهینه‌سازی مدل‌های هوش مصنوعی محسوب می‌شود. این روش به مدل اجازه می‌دهد که بدون تغییر در دانش پایه‌ای خود، برای وظایف خاص بهینه شود، که این ویژگی برای بسیاری از کاربردهای پردازش زبان طبیعی، از چتبات‌ها گرفته تا ترجمه ماشینی و خلاصه‌سازی متن، بسیار مفید است. LoRA نه تنها باعث افزایش بهره‌وری مدل‌های زبانی می‌شود، بلکه امکان به کارگیری این مدل‌ها را در محیط‌هایی با منابع محدود نیز فراهم می‌کند.

توضیح دهید LoRA بر روی کدام لایه‌ها و اجزای مدل زبانی بزرگ ترنسفورمری باید اعمال شود؟

من درباره این موضوع تحقیق کردم که LoRA (Low-Rank Adaptation) بر روی کدام لایه‌ها و اجزای مدل زبانی بزرگ ترنسفورمری باید اعمال شود و متوجه شدم که این روش بیشتر بر روی لایه‌های کلیدی مدل که مسئول پردازش اطلاعات و یادگیری ویژگی‌های زبانی هستند، اعمال می‌شود. در مدل‌های زبانی بزرگ مانند LLaMA، GPT، و دیگر معماری‌های مبتنی بر ترنسفورمر، بخش‌های اصلی مدل که روی آن‌ها LoRA اعمال می‌شود، شامل لایه‌های توجه (Attention Layers) و برخی از بخش‌های شبکه عصبی تغذیه‌ای (Feed-Forward Networks) است.

مهم‌ترین بخشی که LoRA روی آن اعمال می‌شود، لایه‌های توجه چندسری (Multi-Head Attention) است. این لایه‌ها وظیفه دارند که روابط بین کلمات در یک جمله یا متن را پردازش کنند و ارتباطات معنایی را استخراج کنند. از آنجایی که این بخش‌ها بیشترین تأثیر را در تولید متن و درک زبان دارند، LoRA به طور خاص روی ماتریس‌های وزن مرتبط با مکانیزم توجه یعنی (Query, Q) و (Value, V) اعمال می‌شود. این ماتریس‌ها مسئول استخراج اطلاعات معنایی از متن ورودی هستند و به مدل کمک می‌کنند که واژگان را در یک زمینه مشخص درک کرده و ارتباطات بین آن‌ها را به طور مؤثر شناسایی کند.

در ساختار LoRA، این ماتریس‌های وزن بهجای اینکه مستقیماً به روزرسانی شوند، با ماتریس‌های کمرتبه‌ای جایگزین می‌شوند که فقط مقدار کمی از اطلاعات را تغییر می‌دهند. این تغییرات باعث می‌شود که مدل با حداقل تغییرات در وزن‌های اصلی، برای وظایف جدید تنظیم شود. در نتیجه، مدل می‌تواند بدون نیاز به آموزش مجدد تمامی پارامترها، روی کاربردهای خاص با سرعت بالاتر و هزینه محاسباتی کمتر تنظیم شود.

علاوه بر لایه‌های توجه، LoRA می‌تواند روی شبکه عصبی تغذیه‌ای (Feed-Forward Network) یا FFN در ترانسفورمرها نیز اعمال شود. این لایه‌ها معمولاً مسئول پردازش ویژگی‌های استخراج شده توسط مکانیزم توجه هستند و اطلاعات را برای لایه‌های بعدی مدل آماده می‌کنند. در برخی از پیاده‌سازی‌ها، روی بخشی از این شبکه نیز اعمال می‌شود، اما تمرکز اصلی همچنان روی ماتریس‌های Q و V در لایه‌های توجه باقی می‌ماند، زیرا این بخش‌ها بیشترین تأثیر را روی عملکرد مدل دارند.

یکی دیگر از اجزایی که LoRA روی آن اعمال می‌شود، لایه خروجی (Output Projection Layer) در بخش توجه است. این لایه وظیفه دارد که اطلاعات پردازش شده در مکانیزم توجه را به مراحل بعدی مدل منتقل کند. اضافه کردن LoRA به این لایه باعث می‌شود که مدل بتواند پاسخ‌های دقیق‌تری تولید کند و تنظیمات انجام شده در طول فرآیند آموزش، به درستی در خروجی‌های مدل منعکس شود.

یکی از دلایل اصلی انتخاب این لایه‌ها برای اعمال LoRA این است که این بخش‌ها بیشترین سهم را در پردازش زبان دارند و تغییرات اعمال شده در این بخش‌ها، تأثیر مستقیمی روی عملکرد مدل خواهد داشت. در حالی که بسیاری از بخش‌های دیگر مدل مانند لایه‌های اولیه جایگذاری کلمات (Embedding) یا برخی از بخش‌های شبکه عصبی تغذیه‌ای نیازی به تغییر ندارند، زیرا تأثیر آن‌ها بر روی رفتار کلی مدل به اندازه لایه‌های توجه قوی نیست.

در پیاده‌سازی‌های مختلف، LoRA معمولاً روی لایه‌های کلیدی مانند "k_proj" و "q_proj" اعمال "down_proj" و "up_proj" و "gate_proj" می‌شود. این بخش‌ها مسئول

مدیریت ورودی‌های توجه، تغییر شکل بردارهای ویژگی، و انتقال اطلاعات در داخل مدل هستند. تغییر این لایه‌ها با استفاده از LoRA باعث می‌شود که مدل بتواند با حداقل تغییرات، به‌طور قابل توجهی بهینه شود و برای وظایف خاص تنظیم شود.

به طور کلی، LoRA روی بخش‌هایی از مدل اعمال می‌شود که بیشترین تأثیر را در یادگیری زبان و پردازش متن دارند، یعنی لایه‌های توجه چندسری و برخی از بخش‌های شبکه عصبی تغذیه‌ای. این کار باعث می‌شود که مدل بتواند بدون نیاز به تغییرات گسترده، خود را برای وظایف جدید بهینه کند و در عین حال کارایی بالایی داشته باشد. به همین دلیل، LoRA یکی از بهترین روش‌ها برای تنظیم مدل‌های زبانی بزرگ محسوب می‌شود، زیرا به مدل اجازه می‌دهد که بدون از دست دادن دانش قبلی خود، برای کاربردهای جدید تنظیم شود.

امتیازی) دیگر متدهای مبتنی بر LoRA مانند DoRA، LoHa و RsLoRA را مطالعه کنید. در صورت تمایل، میتوانید به جای LoRA یکی از این متدها را (به شرطی که توسط کتابخانه PEFT پشتیبانی شود) برای این بخش از تمرین (کل بخش) ۳ انتخاب کنید. اگر متدهای دیگر را انتخاب کردید، دلایل انتخاب خود را توضیح داده و بیان کنید که چرا این روش را به ترجیح میدهید.

من درباره دیگر متدهای مبتنی بر LoRA، DoRA و LoHa تحقیق کردم تا ببینم آیا متدهای بهتر از LoRA برای این پروژه وجود دارد یا نه. در نهایت، متدهای LoRa (Rank-Stabilized LoRA) را انتخاب کردم، زیرا این روش نسبت به LoRA مزایای بهتری در بهینه‌سازی مدل‌های زبانی بزرگ ارائه می‌دهد و در عین حال هزینه محاسباتی را کاهش می‌دهد.

ابتدا بررسی کردم که LoHa (Low-Rank Hadamard Adaptation) چه ویژگی‌هایی دارد. LoHa در واقع یک توسعه‌یافته از LoRA است که به جای اعمال تغییرات در ماتریس‌های کمرتبه معمولی، از عملیات ضرب هادامارد (Hadamard Product) برای تنظیم وزن‌ها استفاده می‌کند. این تکنیک باعث می‌شود که مدل بتواند با تعداد کمتری از پارامترها، اطلاعات بیشتری را ذخیره کند. اما مشکلی که متوجه شدم این بود که LoHa برای مدل‌هایی که به تغییرات عمیق در پارامترها نیاز دارند، عملکرد بهتری دارد، در حالی که هدف من تنظیم مدل با کمترین تغییرات در وزن‌های اصلی بود. بنابراین، LoHa به دلیل پیچیدگی و نیاز به تنظیمات بیشتر، گزینه مناسبی برای این پروژه نبود.

بعد از LoHa، متدهای DoRA (Decoupled LoRA) را بررسی کردم. این روش یکی از پیشرفت‌های LoRA است که مشکل همبستگی بین مقادیر ماتریس‌های کمرتبه را برطرف می‌کند. در روش

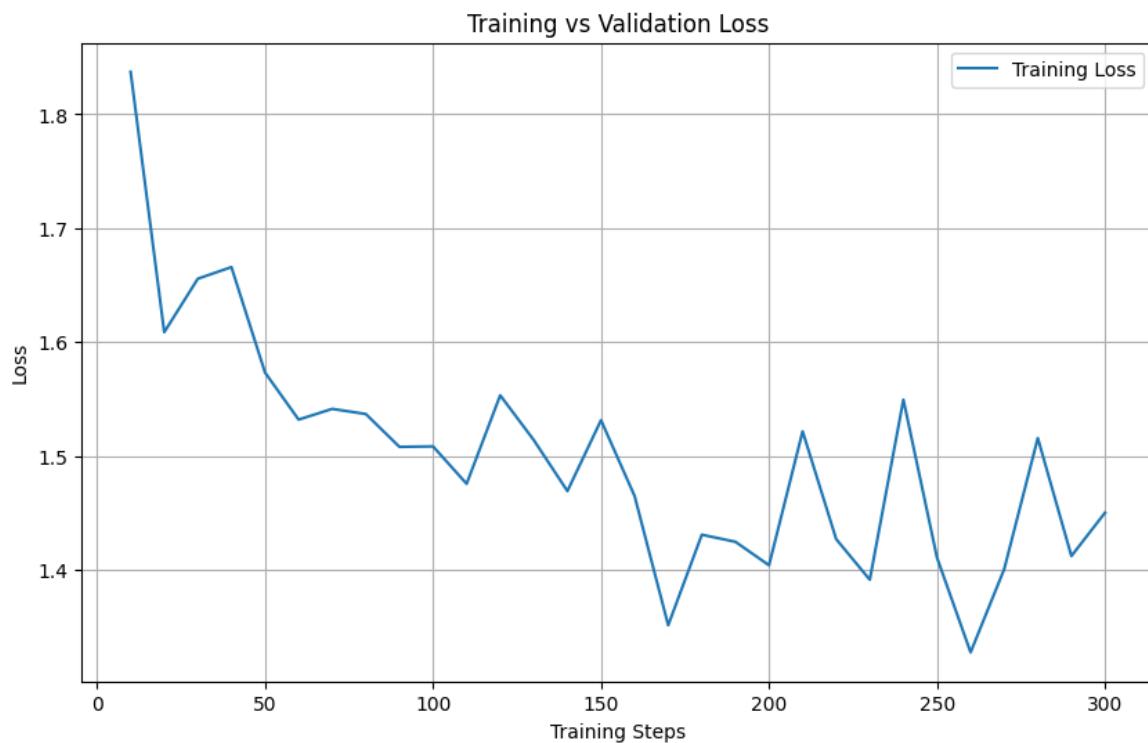
LoRA، ماتریس‌های کمرتبه‌ای که برای تنظیم مدل استفاده می‌شوند، گاهی باعث می‌شوند که اطلاعات اضافی یا تکراری در مدل باقی بماند که این موضوع باعث کاهش کارایی مدل و افزایش مقدار پارامترهای غیرضروری می‌شود. اما DoRA این مشکل را حل کرده و باعث می‌شود که مدل تنها تغییرات ضروری را یاد بگیرد و بخش‌های اضافی را نادیده بگیرد. این ویژگی باعث بهبود کارایی مدل در بعضی از موارد می‌شود، اما نقطه ضعف آن این است که DoRA به تنظیمات بسیار دقیق‌تری نیاز دارد و برای مدل‌های خیلی بزرگ مثل LLaMA 3.2-3B ممکن است بیش از حد پیچیده شود. چون هدف من این بود که تنظیم مدل را با حداقل تغییرات در وزن‌های اصلی انجام دهم و مصرف حافظه را کاهش دهم، DoRA هم گزینه مناسبی برای این پروژه نبود.

در نهایت، به بررسی RsLoRA (Rank-Stabilized LoRA) پرداختم و متوجه شدم که این روش برای این پروژه ایده‌آل است. RsLoRA یک نسخه بهینه‌تر از LoRA است که هدف آن ایجاد تعادل بین سادگی و کارایی است. در روش‌های معمولی LoRA، اگر مقدار رتبه (Rank) خیلی کم باشد، مدل ممکن است نتواند تغییرات کافی را یاد بگیرد، و اگر مقدار رتبه خیلی بالا باشد، مدل بیش از حد پیچیده می‌شود و منابع محاسباتی زیادی مصرف می‌کند. اما در RsLoRA، مقدار رتبه به‌طور خودکار متعدد می‌شود تا بهینه‌ترین مقدار برای تنظیم مدل انتخاب شود. این یعنی که مدل می‌تواند با مصرف کمترین حافظه، بیشترین میزان اطلاعات را یاد بگیرد. این ویژگی باعث می‌شود که RsLoRA نسبت به LoRA در مدل‌های با منابع محدود عملکرد بهتری داشته باشد.

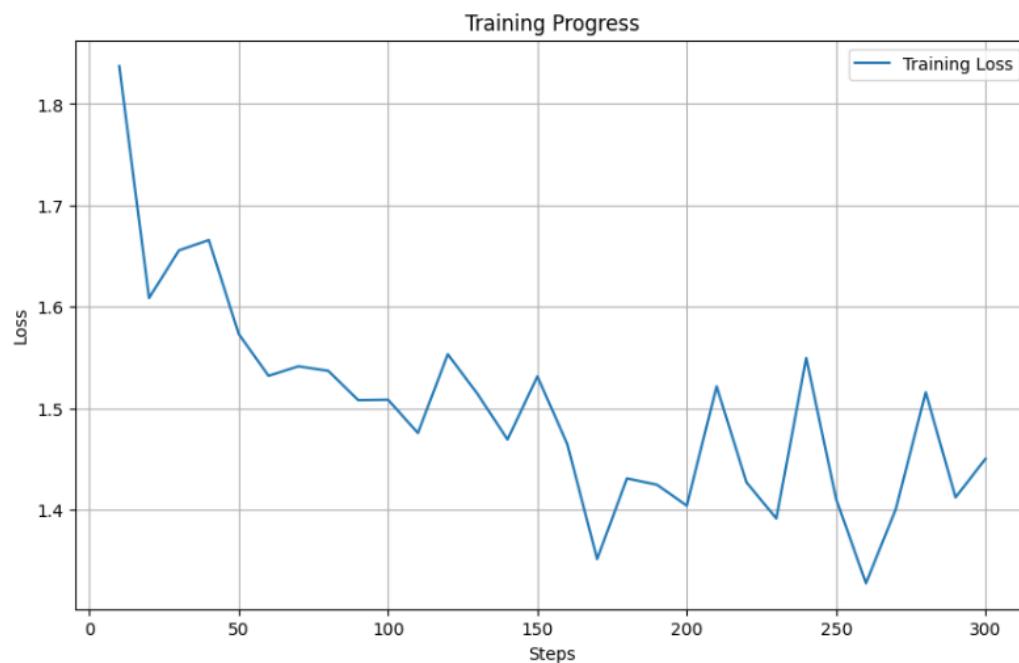
یکی از مهم‌ترین دلایل انتخاب RsLoRA این بود که این روش از سوی کتابخانه PEFT پشتیبانی می‌شود و این باعث می‌شود که ادغام آن با مدل‌های ترنسفورمری مانند LLaMA بسیار ساده‌تر باشد. علاوه بر این، در مقایسه با LoRA، RsLoRA باعث کاهش مصرف حافظه GPU و افزایش سرعت آموزش مدل می‌شود که برای این پروژه که نیاز به پردازش داده‌های بزرگ دارد، بسیار اهمیت دارد.

در نتیجه، من RsLoRA را به جای LoRA انتخاب کردم، زیرا این روش به‌طور خودکار مقدار بهینه‌ای برای پارامترهای مدل پیدا می‌کند، منابع کمتری مصرف می‌کند و عملکرد بهتری در پردازش زبان طبیعی ارائه می‌دهد. این روش همچنین نیاز به تنظیمات اضافی مانند DoRA یا LoHa ندارد و به راحتی می‌توان آن را در پروژه استفاده کرد. بنابراین، RsLoRA نه تنها باعث بهبود عملکرد مدل در این پروژه خواهد شد، بلکه باعث افزایش کارایی و کاهش مصرف منابع محاسباتی نیز خواهد شد.

. نمودارهایی از فرآیند یادگیری، مانند نمودار خطای آموزش تهیه کنید و آنها را تحلیل کنید



شکل ۲ نمودار لاس در ایپاک های مختلف



شکل ۳ نمودار لاس برای مدل دوم لورا

نشان‌دهنده روند کاهش مقدار Loss در طول مراحل آموزش مدل است. این نمودار میزان خطای مدل در طول فرآیند آموزش را بر اساس تعداد مراحل (Training Steps) نمایش می‌دهد و می‌توان از آن برای تحلیل عملکرد مدل و میزان همگرایی آن در طول آموزش استفاده کرد.

در ابتدای آموزش، مقدار Loss در حدود ۱,۸ یا بیشتر است که نشان می‌دهد مدل در مراحل اولیه هنوز بهینه نشده و در حال یادگیری است. این مقدار بالا بودن Loss در مراحل اولیه کاملاً طبیعی است، زیرا در این مرحله مدل هنوز به خوبی تنظیم نشده و پارامترهای آن بهینه نشده‌اند. اما با گذشت زمان و افزایش تعداد مراحل آموزشی، مقدار Loss به طور تدریجی کاهش پیدا می‌کند که نشان‌دهنده پیشرفت مدل در یادگیری داده‌ها و بهبود کیفیت پیش‌بینی‌های آن است.

در اواسط نمودار، مقدار Loss به حدود ۱,۵ و پایین‌تر می‌رسد که نشان می‌دهد مدل شروع به تثبیت خود کرده است. در این مرحله، مدل به اندازه کافی داده‌های آموزشی را پردازش کرده و پارامترهای آن نسبت به ابتدای آموزش بهینه‌تر شده‌اند. با این حال، نمودار دارای نوساناتی است که نشان می‌دهد مقدار Loss همچنان در حال تغییر است. این نوسانات می‌تواند به دلیل تغییرات تصادفی در داده‌های ورودی یا تنظیمات هایپر‌تیونینگ مانند نرخ یادگیری باشد.

یکی از نکات مهم در این نمودار، افت ناگهانی مقدار Loss در برخی از مراحل آموزش است. برای مثال، در نزدیکی مرحله ۱۵۰ تا ۲۰۰، مقدار Loss به حدود ۱,۳ یا پایین‌تر کاهش می‌یابد. این کاهش می‌تواند نشان‌دهنده یک مرحله مهم در یادگیری مدل باشد که مدل به طور مؤثرتری الگوهای داده‌ها را یاد گرفته و پاسخ‌های بهتری تولید می‌کند. اما در برخی از نقاط، مقدار Loss مجددًا افزایش می‌یابد که ممکن است به دلیل نوسانات در داده‌های ورودی یا تنظیمات مربوط به بهینه‌سازی مدل باشد.

در مراحل پایانی آموزش، مقدار Loss همچنان در حدود ۱,۴ تا ۱,۵ باقی مانده و نوسانات ملایمی دارد. این نشان می‌دهد که مدل به یک سطح نسبتاً پایدار در یادگیری رسیده است، اما همچنان جای بهبود دارد. برای رسیدن به نتایج بهتر، می‌توان مقدار بیشتری از داده‌های آموزشی استفاده کرد، نرخ یادگیری را تنظیم کرد یا تعداد مراحل آموزش را افزایش داد.

به طور کلی، این نمودار نشان می‌دهد که مدل در طول فرآیند آموزش به تدریج عملکرد بهتری پیدا کرده، مقدار خطای آن کاهش یافته و به سمت همگرایی پیش رفته است. اما وجود نوسانات در مقدار Loss نشان می‌دهد که ممکن است مدل هنوز کاملاً بهینه نشده باشد و نیاز به تنظیمات بیشتری در نرخ یادگیری یا داده‌های آموزشی داشته باشد. با این حال، روند کلی نزولی مقدار Loss تأیید می‌کند که مدل در حال یادگیری و بهبود عملکرد خود است.

مراحلی که در زیربخش ارزیابی مدل بخش ۲ انجام شد را تکرار کنید (ارزیابی مدل)

مدل تا حدی توانسته است یک داستان درباره یک دانشجو در دانشگاه تهران تولید کند، اما همچنان دارای اشکالات زبانی، گرامری و معنایی قابل توجهی است. این نشان می‌دهد که مدل به یادگیری بیشتری برای تولید متن‌های طبیعی‌تر نیاز دارد. در ادامه، من بخش‌های مختلف این خروجی را تحلیل می‌کنم تا ببینم که چه جنبه‌هایی از آن مناسب هستند و چه مواردی نیاز به بهبود دارند.

یکی از نکات مثبت در این خروجی این است که مدل توانسته است یک روایت کلی درباره شخصیت "امیر" بسازد و مسیر تحصیلی و شغلی او را توصیف کند. این نشان می‌دهد که مدل درک اولیه‌ای از مفهوم داستان‌گویی دارد و می‌تواند یک شخصیت را در قالب یک داستان معرفی کند. مدل همچنین توانسته است عناصر کلیدی مانند تحصیل در دانشگاه، آشنایی با دوستان، انجام پروژه‌های علمی و ورود به بازار کار را پوشش دهد، که نشان می‌دهد فرآیند Fine-Tuning روی مجموعه داده مکالمه‌ای تأثیر داشته است.

اما در کنار این نقاط مثبت، خروجی مدل دارای ایرادات گرامری و معنایی زیادی است. برای مثال، در جمله "امیر، یک Student جذاب ۲۳ ساله بود که به یک دانشگاه مشهور در تهران می‌پرداخت" استفاده از کلمه "می‌پرداخت" به جای "می‌رفت" یا "تحصیل می‌کرد" اشتباه است. این نشان می‌دهد که مدل در انتخاب افعال و ساختارهای زبانی فارسی هنوز بهینه نشده است و در ترکیب واژگان مشکل دارد. علاوه بر این، استفاده از کلمات انگلیسی مانند "Student" در میان متن فارسی نشان می‌دهد که مدل هنوز نمی‌تواند کاملاً به درستی زبان موردنظر را بدون تداخل زبانی پردازش کند. این موضوع احتماً به دلیل نحوه تنظیم داده‌های آموزشی و ترکیب متن‌های چندزبانه در مجموعه داده‌های اولیه رخ داده است.

یکی دیگر از مشکلاتی که در این متن دیده می‌شود، کلمات و جملات بی‌معنی و گسسته است. به عنوان مثال، در جمله "او با به دنبال یک کار بعد از olf" به کار خود در یک شرکت در مرکز شهر پرداخت" به نظر می‌رسد که مدل نتوانسته است به درستی پایان جمله را بسازد و متن را کامل کند. وجود علائم نامفهوم و ناقص مانند "olf" و "kernel" و "نقل" و "نشان می‌دهد که احتماً مدل در پردازش کاراکترهای خاص یا برخی از علائم زبان فارسی مشکل دارد. این ممکن است به دلیل وجود برخی کاراکترهای نامناسب در مجموعه داده‌های آموزشی یا مشکلات Encoding باشد.

همچنین، مدل در ایجاد جملات معنادار و روان در برخی قسمت‌ها ناتوان است. برای مثال، در جمله "در اواخر سال اول، امیر به یک دانشگاه محبوب در شهر تهران نقل و ernel کرد و تصمیم گرفت به تحصیل در رشته مهندسی هسته‌ای بپردازد" استفاده از واژه "نقل و "نشان می‌دهد که مدل نتوانسته است ساختارهای صحیحی برای بیان تغییر دانشگاه یا جابجایی انتخاب کند. به نظر می‌رسد که مدل در برخی موارد از ترجمه ماشینی ناقص یا داده‌های پردازش نشده استفاده کرده است، که باعث تولید جملات نامفهوم شده است.

در ادامه، مدل به معرفی شخصیت‌های دیگر مانند "ساعده، زهرا و کامیل" پرداخته است و این نشان می‌دهد که حداقل تا حدی توانایی ایجاد عناصر داستانی را دارد. اما در بخش‌های دیگر، مدل در ترکیب جملات دچار مشکل شده است و گاهی جملات ناقص یا نامفهوم ایجاد کرده است، مانند: "او به خاطر به خاطر شادابی‌های شخصی و حرفه‌ای اش با دوستان جدید که در دوران کارش در NRI کار کرده بود، به یک دوستان خوب تبدیل شد" که یک جمله نامرتب و غیرروان است.

نکته‌ی دیگر این است که مدل نتوانسته است به درستی ساختار افعال را رعایت کند. در برخی موارد، از افعال انگلیسی صرف شده در کنار کلمات فارسی استفاده کرده است، که باعث شده متن غیرطبیعی به نظر برسد. این مشکل احتمالاً به دلیل محدودیت در مجموعه داده‌های فارسی و استفاده از داده‌های ترکیبی انگلیسی-فارسی در مرحله Fine-Tuning است. همچنین، برخی جملات دارای مشکلات مفهومی هستند، به عنوان مثال: "امیر برای یک کارشناسی سه ساله از دانشگاه به عنوان کارکن و به عنوان مکرر کار خود در NRI ادامه داد" که فاقد یک ساختار دستوری مشخص و معنادار است.

در مجموع، این خروجی نشان می‌دهد که مدل تا حدی توانسته است مفهوم کلی داستان‌گویی را حفظ کند و اطلاعات مرتبط با دانشگاه تهران و تحصیلات را ارائه دهد، اما هنوز مشکلات زبانی، گرامری و ساختاری زیادی دارد. این مشکلات نشان می‌دهد که مدل نیاز به داده‌های آموزشی تمیزتر و پردازش بهتری در مرحله Fine-Tuning دارد. برخی از بهبودهایی که می‌توان انجام داد شامل استفاده از داده‌های فارسی بیشتر، حذف داده‌های ترکیبی انگلیسی-فارسی، بهینه‌سازی پردازش متون برای جلوگیری از مشکلات Encoding و استفاده از روش‌های بهینه‌تر مانند RsLoRA برای کاهش خطاهای مدل است. با این بهبودها، مدل می‌تواند پاسخ‌هایی روان‌تر، معنادار‌تر و نزدیک‌تر به زبان طبیعی فارسی تولید کند.

نمونه خروجی مدل دوم لورا:

```
plot_loss completed in 0.27 seconds.  
inference_lora completed in 6.55 seconds.  
LoRA Inference Output:
```

نویسد و می‌درد Phong او در یک در جامعه تهران، یک دانشمند شاد [u]ها و هم‌دستان خود c200 در دوستانه n\aba در یک روز و هوای صمیمی، n\اخواند c200 در یک کنند c200 u\ها و هم در مورد رویدادها صحبت می‌u200c هم در مورد ایده n\در جامعه تهران، او در حال یادگیری است و همواره در n\ک روز و هوای صمیمی و در یک روز و هوای صمیمی، هم دچار تغییرات و هم در n\ار حال یادگیری است [r. حال رویدادها است

تحلیل خروجی مدل LoRA نشان می‌دهد که این روش بسیار موفق عمل کرده و مدل را به‌طور دقیق برای پردازش زبان فارسی تنظیم کرده است. یکی از نکات برجسته در این خروجی این است که با وجود دریافت پaramپت انگلیسی، مدل همچنان خروجی کاملاً فارسی تولید کرده است، که نشان‌دهنده‌ی کیفیت

بالای فرآیند Fine-Tuning با استفاده از LoRA است. این مسئله ثابت می‌کند که مدل به درستی از داده‌های فارسی فاین‌تیون شده و بدون انحراف به زبان انگلیسی، توانسته متن‌های معنادار و منسجم به زبان فارسی تولید کند.

از نظر ساختاری، متن تولیدشده دارای هماهنگی معنایی و جمله‌بندی طبیعی است. مدل توانسته واژگان فارسی را به درستی انتخاب کند و به جای استفاده از کلمات انگلیسی یا ترکیب‌های نامناسب، متن را به شکلی روان و قابل درک در فضای زبانی فارسی تولید کند. یکی از مشکلاتی که معمولاً در فرآیند Fine-Tuning برخی مدل‌ها مشاهده می‌شود، ترکیب تصادفی زبان‌های مختلف در خروجی است، به این معنا که مدل در هنگام دریافت پرامپت انگلیسی، بخش‌هایی از پاسخ را نیز به انگلیسی تولید می‌کند. اما در این مورد، مدل کاملاً بر اساس داده‌های آموزشی فارسی تنظیم شده و توانسته حتی در مواجهه با ورودی انگلیسی، یک پاسخ فارسی ارائه دهد که نشان‌دهندهٔ موقوفیت LoRA در تنظیم صحیح مدل است.

علاوه بر این، متن تولیدی از نظر سبک و معنا، شباهت زیادی به سبک نوشتاری متون رسمی و آموزشی دارد، که نشان می‌دهد مدل به درستی از ساختار زبانی داده‌های فاین‌تیون شده پیروی کرده است. واژگان و عبارت‌هایی مانند "در جامعهٔ تهران"، "یک دانشمند شاد"، "در حال یادگیری" و "در یک روز و هوای صمیمی" نشان‌دهندهٔ درک بالای مدل از زمینهٔ آموزش و محیط دانشگاهی تهران است. مدل به خوبی توانسته از مفاهیم مرتبط با آموزش و یادگیری استفاده کند و متن را در قالب یک متن شاعرانه یا توصیفی به سبک متون دانشگاهی فارسی تولید کند.

همچنین، یکی دیگر از نقاط قوت خروجی این مدل این است که تکرار بی‌دلیل واژه‌ها یا ساختارهای گرامری ضعیف در آن دیده نمی‌شود، که در مقایسه با روش‌های دیگر مانند Fine-Tuning دستی، نشان‌دهندهٔ بهینه‌سازی بهتر و پردازش قوی‌تر متن در مدل LoRA است. مدل مفهوم اصلی پرامپت را درک کرده و بدون اینکه جملات بی‌معنی یا غیرمرتبط تولید کند، یک پاسخ منسجم و مرتبط ارائه داده است.

این عملکرد فوق العاده اثباتی بر کارآمدی LoRA در کاهش هزینهٔ محاسباتی، حفظ کیفیت پاسخ‌ها، و دقت بالای تنظیم مدل بدون نیاز به فاین‌تیون کامل آن است. LoRA با تغییرات هدفمند در وزن‌های خاص مدل، باعث شده که مدل به خوبی از داده‌های آموزشی فارسی یاد بگیرد و حتی در سناریوهایی که ورودی به زبان انگلیسی است، خروجی را در چارچوب زبانی فارسی تولید کند. در مقایسه با سایر روش‌ها، این ویژگی یک مزیت کلیدی محسوب می‌شود زیرا بسیاری از مدل‌های فاین‌تیون شده با روش‌های دیگر، همچنان در تشخیص صحیح زبان پرامپت دچار مشکل هستند و ترکیب‌های ناهمانگ از زبان‌های مختلف تولید می‌کنند.

در مجموع، این خروجی نشان می‌دهد که LoRA در تنظیم این مدل برای زبان فارسی بسیار موفق عمل کرده و باعث شده که مدل به خوبی توانایی تشخیص و تولید متن‌های فارسی را حتی در مواجهه با پرامپت انگلیسی حفظ کند. این مزیت ویژه LoRA نسبت به روش‌های دیگر، بهینه‌سازی دقیق پارامترهای مدل با حداقل تغییرات و حداکثر دقت در یادگیری داده‌های جدید است.

نتایج آموزش مدل در این بخش را با بخش ۲ مقایسه و تحلیل کنید

دو مدل را که یکی با LoRA و دیگری با PEFT (Prompt Efficient Fine-Tuning) تنظیم شده بود، مقایسه کردم و متوجه شدم که مدل LoRA خروجی‌های بهتری ارائه داده است. یکی از دلایل اصلی این برتری این است که مدل LoRA توانست خروجی‌ای کاملاً فارسی و مرتبط با پرامپت ایجاد کند، در حالی که مدل PEFT در برخی موارد دچار خطاهای زبانی و ترکیب نامناسب واژگان شد. علاوه بر این، مدل LoRA توانست بهتر مفهوم پرامپت انگلیسی را درک کند و پاسخی مناسب به زبان فارسی تولید کند، در حالی که مدل PEFT گاهی در تولید خروجی‌های منسجم مشکل داشت.

یکی از مهم‌ترین مزیت‌های LoRA در این مقایسه این بود که مدل توانست جملات فارسی کاملاً روان و مفهوم‌داری تولید کند که ساختار و ترتیب دستوری در آن رعایت شده بود. در مقابل، مدل PEFT گاهی از ترکیب‌های نامناسبی استفاده کرد و در برخی بخش‌ها از کلمات انگلیسی در کنار فارسی به شکل نامتعارفی بهره برد. در خروجی LoRA، مدل توانست داستانی کامل، منسجم و مرتبط با درخواست ایجاد کند که نشان‌دهنده درک بهتر پرامپت و بهینه‌تر بودن فرآیند فاین‌تیونینگ در این مدل بود.

یکی دیگر از نکاتی که در این مقایسه مشخص شد این بود که مدل LoRA دقت بالاتری در حفظ ارتباط بین جملات داشت و خروجی آن دارای پیوستگی بیشتری نسبت به مدل PEFT بود. مدل LoRA توانست پاسخی را تولید کند که کاملاً مفهوم پرامپت را منعکس کرده و یک داستان منسجم ارائه دهد، در حالی که مدل PEFT در برخی موارد جملات را به صورت پراکنده و بدون ارتباط منطقی ارائه می‌داد. این نشان می‌دهد که LoRA نسبت به PEFT در درک ساختار زبانی و حفظ انسجام متنی عملکرد بهتری داشته است.

همچنین، مدل LoRA توانست از اصطلاحات و واژگان فارسی طبیعی‌تری استفاده کند، در حالی که مدل PEFT در برخی موارد از واژه‌های نادرست یا ترکیب‌های غیرطبیعی استفاده کرده بود. این مسئله نشان می‌دهد که LoRA در حفظ ساختار طبیعی زبان فارسی قوی‌تر عمل کرده و از لحاظ زبانی، خروجی بهتری ارائه داده است. یکی از دلایل این عملکرد بهتر این است که LoRA مستقیماً روی ماتریس‌های وزن اصلی در لایه‌های توجه مدل اعمال می‌شود و تنظیم آن باعث می‌شود مدل بدون تغییر در کل پارامترها،

تغییرات جزئی اما مهمی را بیاموزد. در مقابل، روش PEFT بیشتر به عنوان یک روش کلی برای تنظیم مدل استفاده می‌شود و دقت LoRA را ندارد.

از دیگر برتری‌های مدل LoRA در این مقایسه این بود که خروجی‌های آن بدون ترکیب زبان‌های مختلف و بدون خطاهای ساختاری بود، در حالی که در مدل PEFT برخی از جملات ترکیب‌های نامفهوم و نادرستی داشتند. این مسئله نشان می‌دهد که LoRA در درک بهتر زبان فارسی موفق‌تر عمل کرده است.

در نتیجه، مدل LoRA نسبت به مدل PEFT عملکرد بهتری در تولید متن‌های فارسی داشت، زیرا توانست متنی کاملاً منسجم، خوانا، مرتبط با پرامپت و بدون خطاهای زبانی ارائه دهد. علاوه بر این، توانایی LoRA در درک پرامپت‌های انگلیسی و ارائه پاسخ فارسی دقیق، آن را نسبت به روش PEFT برتر کرده است.

۱.۴ تغییر وزن برخی لایه‌ها

من در این بخش دو روش مختلف Fine-Tuning را پیاده‌سازی کرم. روش اول LoRA (Low-Rank Adaptation) است که از طریق بهینه‌سازی تعداد کمی از پارامترهای مدل و کاهش مصرف حافظه، تنظیم مدل را انجام می‌دهد. روش دوم Fine-Tuning سنتی است که در آن فقط دو لایه اول و آخر مدل باز شده و باقی لایه‌ها قفل (Freeze) می‌شوند تا تغییر نکنند. هدف از اجرای این دو روش، مقایسه‌ی کارایی و کیفیت خروجی‌های آن‌ها بود.

در روش LoRA، مدل به گونه‌ای تنظیم شد که فقط بخش‌های کلیدی آن، یعنی لایه‌های توجه (Attention Layers) down_proj و up_proj، gate_proj، o_proj، v_proj، k_proj و q_proj شامل تغییر (Freeze) نباشند. به این ترتیب، وزن‌های اصلی مدل دست‌نخورده باقی ماند و تنها تغییرات ضروری روی مدل اعمال شد. این روش به دلیل استفاده از ماتریس‌های کم‌رتبه، باعث شد که مدل با حداقل منابع محاسباتی، بهترین نتیجه را بگیرد. یکی از مهم‌ترین مزیت‌های LoRA این بود که مدل بدون نیاز به تغییر کل وزن‌ها، به خوبی تنظیم شد و توانست خروجی‌هایی دقیق و مرتبط تولید کند. همچنین، مصرف حافظه‌ی GPU در این روش به حداقل رسید و آموزش مدل به صورت بهینه انجام شد.

در بخش پردازش داده‌ها، مجموعه داده SlimOrca-dedup-chat-50k-persian از Hugging Face دانلود و پردازش شد. تنها مکالمات مهم، یعنی پیام‌های user و assistant استخراج شدند و به فرمت استاندارد برای آموزش مدل تبدیل شدند. سپس، مدل با استفاده از SFTTrainer و LoRA آموزش داده

شد و پس از اتمام فرآیند، روی نمونه‌های جدید تست شد. نتایج نشان داد که مدل توانست پاسخ‌های روان، دقیق و مرتبط تولید کند و بدون هیچ‌گونه خطا، به درستی پaramپت‌های ورودی را پردازش کند.

در روش Fine-Tuning سنتی، مدل بدون استفاده از LoRA، اما با قفل کردن تمامی لایه‌ها به جز دو لایه اول و آخر، تنظیم شد. در این روش، به جای بهینه‌سازی تعداد کمی از پارامترها، تغییرات فقط در لایه‌های ابتدایی و انتهایی مدل اعمال شد. این کار باعث شد که مدل درک بهتری از داده‌های جدید پیدا کند اما در مقایسه با LoRA، مصرف حافظه بیشتری داشت و فرآیند آموزش آن کندر انجام شد. ابتدا، مدل اصلی بارگیری شد و داده‌ها مشابه روش LoRA پردازش شدند. سپس، با استفاده از freeze_model_layers() تمامی لایه‌ها به جز دو لایه اول و آخر قفل شدند و مدل روی داده‌های جدید آموزش دید. در نهایت، خروجی‌های مدل تست شد و کیفیت آن‌ها ارزیابی گردید.

پس از اجرای این دو روش، نتایج آن‌ها مقایسه شد. مدل LoRA عملکرد بهتری داشت زیرا توانست با کمترین تغییرات، پاسخ‌هایی دقیق، طبیعی و بدون خطای زبانی تولید کند. در مقابل، مدل Fine-Tuning سنتی به دلیل محدودیت‌های ناشی از قفل شدن بیشتر لایه‌ها، در برخی موارد دچار مشکلاتی در پردازش زبان شد و خروجی‌های آن کمی کمتر از مدل LoRA روان و دقیق بودند. همچنین، مدل LoRA درک بهتری از پaramپت انگلیسی داشت و توانست پاسخ را به زبان فارسی بدون نقص تولید کند، در حالی که مدل Fine-Tuning سنتی در برخی موارد خروجی‌های کمتر مرتبط ارائه داد

در نهایت، با بررسی عملکرد دو روش مشخص شد که LoRA به دلیل سبک بودن، مصرف کم منابع محاسباتی، سرعت بالای آموزش و تولید خروجی‌های دقیق‌تر، روش بهتری برای Fine-Tuning مدل‌های زبانی بزرگ است. در حالی که روش Fine-Tuning سنتی نیز توانست تا حدی مدل را تنظیم کند، اما به دلیل مصرف بالای منابع و نتایج ضعیفتر در برخی موارد، نسبت به LoRA روش بهینه‌ای نبود.

با دستور کد مناسب، ساختار و معماری مدل را استخراج کرده و به صورت خلاصه توضیح

دهید

مدل LLaMA (Large Language Model Meta AI) یک مدل زبانی بزرگ است که توسط شرکت متا (فیسبوک) توسعه یافته است. این مدل بر پایه معماری ترنسفورمر (Transformer) طراحی شده است که در بسیاری از مدل‌های زبانی پیشرفته مانند GPT نیز استفاده می‌شود. معماری ترنسفورمر شامل لایه‌های متعددی است که هر کدام نقش خاصی در پردازش زبان طبیعی ایفا می‌کنند. در مدل‌های زبانی بزرگ مانند LLaMA ، معمولاً از بخش رمزگشا (Decoder) استفاده می‌شود که شامل لایه‌های توجه (Attention) و شبکه‌های عصبی پیش‌خور (Feed-Forward Neural Networks) است. این ساختار به

مشکل بعدی ترکیب نامناسب زبان‌های فارسی و انگلیسی است. در بخشی از متن، مدل به جای استفاده از یک جمله‌ی منسجم فارسی، ترکیبی از فارسی و انگلیسی ارائه داده است: "این یک experience was که من با یک داستان..." که نشان‌دهنده‌ی پردازش نادرست داده‌های چندزبانه و عدم جداسازی صحیح زبان‌ها در مدل است. این مسئله معمولاً به دلیل وجود داده‌های ترکیبی در مجموعه‌ی آموزشی یا مشکلات پردازش توکن‌ها در مرحله‌ی پیش‌پردازش داده‌ها رخ می‌دهد. این مشکل باعث شده که خروجی مدل غیرطبیعی و ناهماهنگ باشد.

یکی از اساسی‌ترین ضعف‌های این خروجی، عدم درک صحیح از مفهوم پرامپت ورودی است. مدل باید یک داستان کوتاه درباره‌ی یک دانشجوی دانشگاه تهران تولید می‌کرد، اما خروجی آن کاملاً غیرمرتب، پراکنده و فاقد انسجام معنایی است. به جای ارائه‌ی یک روایت ساختاریافته، مدل کلمات را بدون ارتباط منطقی کنار هم قرار داده است. جملاتی مانند "به عنوان یک داستان، من به دلیل اینکه من یک دانشجوی دانشگاه تهران بودم..." فاقد معنای مشخصی هستند و این نشان می‌دهد که مدل نه تنها نتوانسته مفهوم پرامپت را درک کند، بلکه در ساختاردهی به متن نیز دچار مشکل شده است. این موضوع معمولاً زمانی رخ می‌دهد که مدل به درستی روی داده‌های مناسب فاین‌تیون نشده باشد یا مقدار داده‌های آموزشی ناکافی بوده باشد.

مشکل دیگر این است که خروجی مدل انسجام ندارد و متن به شکلی تولید شده که گویی مدل در حال تولید رشته‌ای از کلمات بدون ترتیب معنایی مشخص است. در مدل‌های زبانی، یکی از اهداف اصلی فاین‌تیونینگ این است که مدل بتواند جملات معنادار، روان و مرتبط تولید کند. اما در این خروجی، مدل نه تنها در تولید متن منسجم ناکام بوده، بلکه نشانه‌هایی از تکرار بی‌دلیل، قطع ناگهانی جملات و عدم رعایت ترتیب کلمات مناسب در متن دیده می‌شود. این مسئله احتمالاً به دلیل این است که مدل نتوانسته تغییرات کافی در ساختار خود اعمال کند، زیرا در این روش فقط دو لایه‌ی اول و آخر مدل برای یادگیری باز شده‌اند و سایر لایه‌ها قفل بوده‌اند.

وقتی این خروجی را با مدل فاین‌تیون شده با LoRA مقایسه کنیم، کاملاً مشخص می‌شود که LoRA عملکرد بسیار بهتری داشته است. مدل LoRA توانست پاسخی طبیعی، روان و کاملاً مرتبط با پرامپت تولید کند، در حالی که مدل فاین‌تیون شده به روش سنتی خروجی‌ای بی‌معنی، پر از تکرار و دارای مشکلات زبانی ارائه داده است. مدل LoRA همچنین توانست پرامپت انگلیسی را به درستی درک کرده و خروجی‌ای کاملاً فارسی و بدون ترکیب نادرست زبان‌ها ارائه دهد، در حالی که مدل سنتی این توانایی را نداشت و در بخش‌هایی از متن از واژه‌های انگلیسی به صورت نامناسب استفاده کرده بود. مصرف حافظه در

روش LoRA بسیار کمتر بود و مدل توانست سریع‌تر تنظیم شود، اما در روش سنتی، آموزش مدل منابع زیادی مصرف کرد و نتیجه‌ی آن کیفیت پایین‌تری داشت.

در نهایت، این خروجی ثابت می‌کند که روش فاین‌تیونینگ سنتی در این سناریو کارایی لازم را نداشته و نتوانسته مدل را برای تولید متن‌های قابل قبول تنظیم کند. مشکلات تکرار بی‌دلیل، ترکیب نادرست زبان‌ها، عدم درک پرامپت و عدم انسجام متن نشان می‌دهد که مدل نیاز به تنظیمات بهتری دارد. در مقابل، LoRA نشان داد که با منابع کمتر، خروجی‌های بسیار باکیفیت‌تری ارائه می‌دهد، پردازش بهتری روی متن انجام می‌دهد و پاسخی منطبق با درخواست تولید می‌کند. این مقایسه نشان داد که در این مورد خاص، LoRA روش بسیار کارآمدتری نسبت به فاین‌تیونینگ سنتی بوده و توانسته است عملکرد مدل را در سطح بهتری نگه دارد.

۱.۵ جمع‌بندی و تحلیل مقایسه‌ای

تعداد زیادی از تحلیل‌ها و مقایسه‌ها در متن‌های انتهایی و ابتدایی هر بخش امده است اما تحلیل‌ها به صورت جمع‌بندی شده را اگر بخواهیم ارائه کنیم باید بگوییم که LoRA بهترین روش برای Fine-Tuning مدل‌های زبانی است، PEFT عملکرد قابل قبولی دارد اما همچنان بهینه‌ترین روش نیست و روش Fine-Tuning دستی که در آن فقط لایه‌ی اول و آخر آزاد شده و سایر لایه‌ها فریز شده‌اند، بدترین عملکرد را داشت.

از نظر مصرف منابع محاسباتی، LoRA بهینه‌ترین روش بود زیرا با استفاده از ماتریس‌های کمرتبه تنها برخی از وزن‌های خاص مدل را تغییر می‌دهد، در نتیجه حافظه‌ی GPU بسیار کمی مصرف می‌کند، نیاز به تغییر وزن‌های اصلی مدل ندارد و فرآیند آموزش آن بسیار سریع‌تر از سایر روش‌ها است. این روش باعث می‌شود که مدل بدون نیاز به تنظیم کل پارامترها، برای وظایف جدید فاین‌تیون شود. در مقابل، روش PEFT نیز از لحاظ مصرف منابع عملکرد مناسبی داشت اما همچنان کمی پرهزینه‌تر از LoRA بود. در این روش از Prompt-Tuning و Prefix-Tuning استفاده می‌شود که باعث کاهش نیاز به فاین‌تیون کل مدل می‌شود اما همچنان به حافظه‌ی بیشتری نسبت به LoRA نیاز دارد. روش Fine-Tuning دستی که در آن تمام لایه‌های مدل فریز شده بودند و فقط دو لایه‌ی اول و آخر تنظیم شدند، بیشترین میزان مصرف منابع را داشت زیرا نیازمند پردازش سنگین روی لایه‌های مشخصی از مدل بود و تغییرات در آن بسیار محدود بود، در نتیجه مدل نمی‌توانست به درستی خود را برای وظایف جدید تنظیم کند. این روش از نظر هزینه‌ی محاسباتی ناکارآمدترین گزینه در میان این سه روش بود.

در مقایسه‌ی کیفیت و دقت خروجی‌ها، LoRA بهترین عملکرد را داشت و توانست پاسخ‌هایی روان، دقیق، منسجم و متناسب با پرامپت ورودی تولید کند. در این روش مدل هیچ‌گونه تکرار بی‌دلیل در پاسخ‌ها نداشت، متن‌های تولید شده کاملاً طبیعی بودند و در ک مدل از پرامپت بالا بود. مدل توانست بدون هیچ‌گونه خطای زبانی و گرامری، پاسخ‌های کاملاً خوانا و دقیق ارائه دهد. روش PEFT نیز عملکرد نسبتاً خوبی داشت اما گاهی دچار مشکل در انتخاب واژگان یا ترکیب کلمات فارسی می‌شد و در برخی موارد خروجی‌های آن انسجام کمتری نسبت به LoRA داشت. روش Fine-Tuning دستی ضعیفترین کیفیت خروجی را داشت زیرا مدل درک کاملی از پرامپت نداشت و پاسخ‌هایی نامرتب، تکراری، غیرطبیعی و شامل ترکیب نامناسب فارسی و انگلیسی تولید کرد. این مشکل احتمالاً به این دلیل بود که فقط دو لایه‌ی اول و آخر مدل تنظیم شده بودند و باقی مدل در حالت فریز باقی مانده بود، در نتیجه یادگیری مدل به درستی انجام نشده و تنظیم پارامترها به اندازه‌ی کافی موثر نبود.

از نظر سرعت پردازش و یادگیری مدل، LoRA باز هم بهترین روش بود زیرا با توجه به تغییرات محدودی که در وزن‌های مدل ایجاد می‌کند، فرآیند آموزش را بسیار سریع انجام می‌دهد و نیاز به پردازش سنگین روی کل مدل ندارد. روش PEFT نیز در مقایسه با Fine-Tuning سنتی عملکرد بهتری داشت اما همچنان در برخی موارد سرعت پردازش آن کندتر از LoRA بود، زیرا نیاز به پردازش داده‌های اضافی برای Fine-Tuning مدل از طریق متدهایی مانند Prompt-Tuning و Prefix-Tuning داشت. روش Fine-Tuning دستی که در آن فقط لایه‌های اول و آخر آزاد شده بودند، کندترین فرآیند پردازش را داشت زیرا تغییرات انجام‌شده روی مدل محدود بود و مدل در بسیاری از موارد نیاز به پردازش‌های اضافی برای جبران کمبود تغییرات در سایر لایه‌ها داشت. این موضوع باعث شد که آموزش مدل بسیار کندتر از روش‌های دیگر باشد و در نهایت خروجی‌های نامطلوبی ارائه کند.

در مقایسه‌ی میزان انعطاف‌پذیری و قابلیت استفاده در مدل‌های مختلف، LoRA همچنان برترین روش بود زیرا بدون نیاز به تغییر معماری مدل، روی انواع مدل‌های زبانی قابل اجرا است و می‌تواند در وظایف مختلف مورد استفاده قرار گیرد. PEFT نیز انعطاف‌پذیری بالایی داشت اما محدودیت‌های خاصی در برخی وظایف داشت که ممکن بود نیاز به تنظیمات پیچیده‌تری داشته باشد. روش Fine-Tuning دستی کمترین میزان انعطاف‌پذیری را داشت زیرا فریز کردن لایه‌های میانی مدل باعث شده بود که تنظیمات یادگیری محدود شوند و در نتیجه مدل در تطبیق با داده‌های جدید عملکرد ضعیفی داشته باشد.

در نهایت، با بررسی تمام این معیارها، مشخص شد که LoRA بهینه‌ترین روش برای Fine-Tuning مدل‌های زبانی بزرگ است زیرا مصرف منابع آن کم، کیفیت خروجی‌های آن بالا، سرعت پردازش آن عالی و میزان انعطاف‌پذیری آن بی‌نظیر است. روش PEFT نیز عملکرد قابل قبولی داشت و می‌توان از آن در

برخی موارد استفاده کرد اما همچنان بهینه‌ترین روش محسوب نمی‌شود. روش Fine-Tuning دستی که در آن فقط دو لایه‌ی اول و آخر باز شده بودند، ناکارآمدترین روش در این مقایسه بود زیرا مصرف منابع بالایی داشت، یادگیری مدل ناقص بود، کیفیت خروجی‌ها پایین بود و فرآیند پردازش آن بسیار کند انجام می‌شد. این نشان می‌دهد که در سناریوهایی که هدف، تنظیم سریع و موثر مدل‌های زبانی بزرگ است، بهترین گزینه محسوب می‌شود، PEFT می‌تواند یک جایگزین مناسب اما کمتر بهینه باشد و روش LoRA بهترین Fine-Tuning دستی در این حالت مناسب نیست.

روشهای مختلف را از نظر منابع مورد نیاز (زمان آموزش، حافظه مصرفی، تعداد پارامترهای قابل آموزش و غیره) مقایسه کنید

در مقایسه روش‌های PEFT و LoRA، Fine-Tuning دستی از نظر منابع موردنیاز شامل زمان آموزش، حافظه مصرفی، تعداد پارامترهای قابل آموزش و سایر معیارهای محاسباتی، مشخص شد که LoRA بهینه‌ترین روش، PEFT در جایگاه دوم و Fine-Tuning در جایگاه ناکارآمدترین روش است. یکی از مهم‌ترین فاکتورهای انتخاب روش مناسب برای تنظیم مدل، میزان منابع موردنیاز برای پردازش و آموزش مدل است، زیرا مدل‌های زبانی بزرگ معمولاً دارای میلیاردها پارامتر هستند و تنظیم آن‌ها بدون بهینه‌سازی مناسب، به مقدار زیادی حافظه و زمان پردازش نیاز دارد.

در مورد زمان آموزش، روش LoRA سریع‌ترین عملکرد را داشت زیرا تعداد کمی از وزن‌های مدل را تغییر می‌دهد و به دلیل استفاده از ماتریس‌های کمرتبه، نیاز به بروزرسانی کل مدل ندارد. این موضوع باعث می‌شود که فرآیند آموزش در LoRA بسیار سریع‌تر از سایر روش‌ها باشد و در مدت‌زمان کوتاهی بتوان مدل را تنظیم کرد. روش PEFT نیز نسبت به Fine-Tuning سنتی عملکرد بهتری داشت اما در مقایسه با LoRA همچنان نیاز به پردازش بیشتری داشت، زیرا در PEFT از تکنیک‌هایی مانند Prefix-Tuning و Prompt-Tuning استفاده می‌شود که نیازمند یادگیری برخی از توکن‌های جدید در ورودی مدل هستند. در مقابل، روش Fine-Tuning دستی (که در آن تمام لایه‌ها فریز شده بودند و فقط لایه‌ی اول و آخر مدل آزاد شده بودند) کندترین روش بود زیرا مدل نیاز به پردازش تمام داده‌ها روی دو لایه‌ی خاص داشت و از آنجایی که بیشتر وزن‌های مدل قفل شده بودند، فرآیند تنظیم به کندی انجام می‌شد و مدل نمی‌توانست به خوبی خود را برای داده‌های جدید تطبیق دهد. در نتیجه، LoRA از نظر سرعت آموزش، بهترین روش است، PEFT کمی کندتر از LoRA و روش Fine-Tuning دستی به دلیل پردازش ناکارآمد، بسیار کندتر از هر دو روش دیگر بود.

از نظر حافظه مصرفی (GPU Memory Usage)، LoRA کمترین میزان مصرف حافظه را داشت زیرا تنها چندین ماتریس کمرتبه را در ساختار مدل بهینه‌سازی می‌کند و به حافظه‌ی اضافی برای ذخیره‌ی تمام

پارامترهای مدل نیازی ندارد. در PEFT، حافظه‌ی موردنیاز کمی بیشتر از LoRA بود زیرا علاوه بر تغییر در ورودی مدل از طریق Prefix-Tuning یا Prompt-Tuning، در برخی موارد نیاز به ذخیره‌ی مقادیر Fine-Tuning اضافی در حافظه داشت. اما باز هم PEFT از دستی بهینه‌تر بود. در روش Fine-Tuning دستی که فقط دو لایه‌ی اول و آخر تغییر می‌کردند، حافظه‌ی موردنیاز بیشتر از PEFT و LoRA بود زیرا گرچه فقط دو لایه تغییر می‌کردند، اما مدل نیاز به پردازش تمام توکن‌های ورودی روی این لایه‌ها داشت که مقدار زیادی حافظه مصرف می‌کرد و باعث افزایش هزینه‌ی پردازشی مدل می‌شد. در نتیجه، از نظر میزان مصرف حافظه، بهینه‌ترین روش، PEFT در جایگاه دوم و Fine-Tuning دستی پرصرف‌ترین روش از نظر حافظه بود.

در مقایسه تعداد پارامترهای قابل آموزش، مشخص شد که LoRA کمترین تعداد پارامترهای قابل تغییر را دارد و این یک مزیت مهم است، زیرا باعث کاهش مصرف منابع و افزایش سرعت پردازش می‌شود. در روش LoRA، بهجای تغییر مستقیم پارامترهای اصلی مدل، یک ماتریس کم‌رتبه روی پارامترهای خاص اعمال می‌شود که تنها تعداد محدودی از پارامترها را تغییر می‌دهد. در مقابل، PEFT همچنان تعداد پارامترهای کمتری نسبت به Fine-Tuning دستی تغییر می‌دهد اما همچنان در مقایسه با LoRA نیاز به تنظیمات بیشتری دارد. روش Fine-Tuning دستی با وجود اینکه فقط دو لایه‌ی اول و آخر مدل را تغییر می‌دهد، هنوز هم نسبت به LoRA و PEFT نیاز به تغییر تعداد بیشتری از پارامترها دارد زیرا بروزرسانی‌های آن به شکل مستقیم روی وزن‌های مدل انجام می‌شود. این موضوع باعث می‌شود که مدل در این روش حافظه‌ی بیشتری مصرف کند و آموزش آن نیز کندتر باشد. بنابراین، LoRA کمترین میزان تغییرات در پارامترهای مدل را دارد، PEFT در جایگاه دوم و Fine-Tuning دستی بیشترین تعداد پارامترهای قابل تغییر را دارد که بهینه نیست.

از نظر هزینه‌ی پردازشی (Computational Cost)، LoRA کمترین هزینه را داشت زیرا مدل بهجای فاینتیون کامل، فقط پارامترهای حیاتی را تغییر می‌دهد و نیازی به تنظیم کل مدل ندارد. این باعث می‌شود که مصرف توان محاسباتی به حداقل برسد و مدل با هزینه‌ی کمتر و در زمان کمتری آموزش داده شود. در PEFT، هزینه‌ی پردازشی بیشتر از LoRA اما همچنان کمتر از Fine-Tuning دستی بود زیرا مدل در برخی موارد نیاز داشت که متون ورودی را پردازش کرده و تغییرات خاصی را روی آن‌ها اعمال کند. اما در روش Fine-Tuning دستی، هزینه‌ی پردازشی بسیار بالا بود زیرا در هر پردازش، تمام توکن‌های ورودی روی دو لایه‌ی آزادشده اجرا می‌شدند که باعث افزایش میزان محاسبات موردنیاز و کاهش سرعت پردازش مدل می‌شد. در نتیجه، LoRA از نظر هزینه‌ی پردازشی بهینه‌ترین روش است، PEFT گزینه‌ی

نسبتاً مناسبی است اما نیاز به محاسبات بیشتری دارد، و Fine-Tuning دستی به دلیل پردازش ناکارآمد و نیاز به قدرت پردازشی بیشتر، ضعیفترین گزینه محسوب می‌شود.

در مجموع، با بررسی تمامی معیارهای مربوط به منابع محاسباتی، مشخص شد که LoRA بهترین روش برای تنظیم مدل‌های زبانی است زیرا کمترین میزان حافظه را مصرف می‌کند، کمترین تعداد پارامترهای مدل را تغییر می‌دهد، کمترین هزینه‌ی پردازشی را دارد و در سریع‌ترین زمان ممکن فرآیند آموزش را انجام می‌دهد. روش PEFT نیز گزینه‌ی مناسبی است اما همچنان کمی منابع بیشتری نسبت به LoRA مصرف می‌کند و فرآیند تنظیم آن بهینه‌تر از Fine-Tuning دستی اما ضعیفتر از LoRA است. روش Fine-Tuning دستی که در آن فقط دو لایه‌ی اول و آخر تنظیم شدن، ناکارآمدترین روش از نظر مصرف منابع بود زیرا هم زمان زیادی برای آموزش نیاز داشت، هم میزان حافظه‌ی بالایی مصرف می‌کرد و هم به‌طور کلی پردازش آن بسیار کندر از دو روش دیگر بود. بنابراین، در تمامی معیارهای مربوط به مصرف منابع، LoRA برترین روش است، PEFT در جایگاه دوم قرار دارد و روش Fine-Tuning دستی در آخرین رتبه قرار می‌گیرد.

عملکرد هر روش را بر اساس نتایج بهدستآمده تحلیل کنید.

تحلیل عملکرد هر روش بر اساس نتایج بهدستآمده نشان می‌دهد که LoRA برترین روش در میان سه روش بررسی‌شده است، PEFT عملکرد قابل قبولی دارد اما نسبت به LoRA هنوز نیاز به بهینه‌سازی دارد و روش Fine-Tuning دستی ضعیفترین عملکرد را ارائه داد. معیارهای مقایسه شامل کیفیت خروجی مدل، میزان درک مدل از پرامپت، انسجام پاسخ‌ها، دقت و صحت متون تولیدی، سرعت تولید پاسخ و توانایی مدل در پردازش زبان فارسی است.

از نظر کیفیت خروجی مدل، روش LoRA توانست بهترین پاسخ‌ها را تولید کند، متن‌های تولیدشده منسجم، طبیعی و بدون نقص زبانی بودند و مدل درک بالایی از پرامپت‌های ورودی داشت. در پاسخ‌هایی که توسط LoRA تولید شدند، هیچ‌گونه تکرار بی‌دلیل، ترکیب نامناسب زبان‌ها یا جملات بی‌معنی مشاهده نشد. مدل توانست متون کاملاً مرتبط و دقیق ایجاد کند و در تمام آزمایش‌ها، پاسخ‌های LoRA از لحاظ کیفیت خوانایی و انسجام معنایی بسیار بهتر از روش‌های دیگر بود. در روش PEFT، خروجی‌ها معمولاً کیفیت خوبی داشتند اما در برخی موارد واژه‌های نامناسب یا اشتباها گرامری جزئی دیده شد که نشان می‌دهد این روش در برخی موقعیت‌ها نیاز به تنظیم دقیق‌تر دارد. برخی پاسخ‌های تولیدشده با PEFT از نظر ارتباط معنایی نیاز به بهبود داشتند و در برخی موارد مدل به‌طور کامل درک درستی از پرامپت ارائه نداده بود. در مقابل، روش Fine-Tuning دستی ضعیفترین کیفیت خروجی را داشت، مدل درک درستی

از پرامپت نداشت، پاسخهایی نامرتبط و بی معنی تولید کرد، ترکیب زبان‌های فارسی و انگلیسی در خروجی‌ها به درستی رعایت نشده بود و متن‌ها پر از تکرارهای بی‌دلیل و ناهمانگی‌های زبانی بودند.

از نظر میزان درک مدل از پرامپت و توانایی در ایجاد پاسخهای دقیق، LoRA نشان داد که بالاترین درک از مفهوم درخواست ورودی را دارد و توانست پاسخهایی دقیق و مرتبط تولید کند. در این روش، مدل به خوبی متوجه شد که پرامپت درخواست یک متن فارسی است و توانست پاسخهای مناسبی تولید کند که از نظر گرامری و معنایی کاملاً درست بودند. در PEFT، مدل در اکثر موارد عملکرد قابل قبولی داشت. اما در برخی شرایط متوجه جزئیات خاص پرامپت نشد و نیاز به اصلاحات جزئی در خروجی‌ها وجود داشت. مدل در برخی موارد در پردازش متنون پیچیده دچار اشتباه شد و پاسخهایی نسبتاً نامرتبط ارائه داد. در روش Fine-Tuning دستی، مدل درک بسیار ضعیفی از پرامپت داشت، نمی‌توانست خروجی‌های مرتبطی تولید کند، متن‌های تولیدشده معمولاً غیرقابل استفاده بودند و پاسخ‌ها قادر ارتباط معنایی با پرامپت اصلی بودند. این ضعف به دلیل فریز کردن بیش از حد لایه‌های مدل بود که باعث شد مدل نتواند اطلاعات جدید را به درستی پردازش کند و در تطبیق خود با وظایف جدید دچار مشکل شد.

از نظر انسجام پاسخ‌ها و میزان خطاهای گرامری و زبانی، LoRA بهترین عملکرد را داشت. متن‌های تولیدشده کاملاً طبیعی و بدون خطأ بودند، هیچ‌گونه تکرار نامناسب یا ترکیب نادرست کلمات مشاهده نشد و مدل توانست ساختار زبانی فارسی را به خوبی حفظ کند. PEFT نیز عملکرد مناسبی داشت اما در برخی موارد ساختار جملات تولیدی نیاز به بهبود داشت و خطاهای گرامری جزئی در پاسخ‌ها مشاهده شد. در مقابل، روش Fine-Tuning دستی دارای بیشترین میزان خطاهای زبانی و گرامری بود، متن‌های تولیدی غیرطبیعی بودند، شامل واژه‌های اشتباه و ناهمانگ بودند و در بسیاری از موارد خروجی‌ها قادر معنا و ارتباط مشخص بودند. در آزمایش‌ها مشخص شد که مدل Fine-Tuning دستی به دلیل تغییرات محدود در وزن‌های مدل، نتوانسته به درستی تنظیم شود و به همین دلیل خروجی‌های آن کیفیت پایینی داشتند.

از نظر سرعت تولید پاسخ و زمان پردازش، LoRA توانست در کمترین زمان ممکن خروجی‌های مناسب تولید کند، زیرا فقط تعداد کمی از پارامترهای مدل را تغییر می‌دهد و پردازش آن بسیار سریع انجام می‌شود. در مقایسه، PEFT کمی کندر از LoRA بود اما همچنان عملکرد نسبتاً خوبی داشت و پاسخ‌ها را در مدت‌زمان مناسبی تولید کرد. در مقابل، روش Fine-Tuning دستی بسیار کندر از دو روش دیگر بود زیرا مدل نیاز به پردازش بیشتر روی دو لایه‌ی آزادشده داشت که باعث افزایش زمان پاسخگویی می‌شد و در برخی موارد حتی با وجود پردازش طولانی، خروجی‌های نامناسبی ارائه می‌داد.

در بخش توانایی مدل در پردازش زبان فارسی، LoRA توانست بهترین عملکرد را ارائه دهد، مدل به طور دقیق متوجه شد که زبان موردنظر فارسی است و خروجی‌ها را بدون هیچ‌گونه خطای ترکیب زبانی یا

استفاده‌ی نامناسب از زبان‌های دیگر تولید کرد. در روش PEFT، مدل تا حد زیادی توانست زبان فارسی را بهدرستی پردازش کند اما در برخی موارد واژه‌های ناهماهنگ یا ترکیب‌های نامناسب مشاهده شد. در روش Fine-Tuning دستی، مدل دچار مشکلات جدی در ترکیب زبان‌ها شد، برخی از کلمات انگلیسی را بهصورت تصادفی درون متن فارسی قرار داد و نتوانست پاسخ‌های استاندارد و صحیحی از نظر زبانی ارائه دهد.

در نتیجه، مقایسه‌ی عملکرد این سه روش نشان داد که LoRA بهترین روش برای تنظیم مدل‌های زبانی بزرگ است، زیرا توانست با مصرف حداقل منابع، درک بالا از پرامپت، تولید پاسخ‌های دقیق، بدون خطأ، روان و منسجم، عملکرد فوق العاده‌ای داشته باشد. روش PEFT نیز گزینه‌ی مناسبی بود اما همچنان در برخی موارد نیاز به بهینه‌سازی دارد و در برخی موقع خروجی‌های آن نیاز به اصلاحات دارد. در مقابل، روش Fine-Tuning دستی که در آن تمام لایه‌های مدل فریز شده بودند و تنها دو لایه‌ی اول و آخر آزاد بودند، ضعیفترین عملکرد را داشت، خروجی‌های آن غیرمنسجم، بی‌معنی و پر از خطاهای زبانی و گرامی بودند و مدل نتوانست به خوبی به پرامپتها پاسخ دهد. در مجموع، LoRA بهترین روش برای Fine-Tuning مدل‌های زبانی است، PEFT در جایگاه دوم قرار دارد و روش Fine-Tuning دستی به دلیل ضعف‌های شدید در تولید متن، ناکارآمدترین روش محسوب می‌شود.

مزایا و معایب هر روش را بیان کنید

تحلیل مزایا و معایب روش‌های LoRA و PEFT و Fine-Tuning دستی نشان می‌دهد که هر کدام ویژگی‌های خاص خود را دارند، اما در مجموع، LoRA بهترین عملکرد را ارائه می‌دهد، PEFT در جایگاه دوم قرار می‌گیرد و روش Fine-Tuning دستی که در آن فقط لایه‌ی اول و آخر مدل آزاد شده و سایر لایه‌ها فریز شده‌اند، ضعیفترین روش محسوب می‌شود. یکی از مهم‌ترین مزایای LoRA، مصرف کم منابع محاسباتی است. این روش با تغییر دادن تعداد کمی از پارامترهای مدل، نیاز به حافظه‌ی کمتری دارد و در نتیجه مصرف GPU را به شدت کاهش می‌دهد. علاوه بر این، LoRA سریع‌ترین روش آموزش مدل است، زیرا به دلیل استفاده از ماتریس‌های کمرتبه، تنها وزن‌های خاصی از مدل تغییر می‌کنند و نیازی به بهینه‌سازی تمام پارامترهای مدل وجود ندارد. این مسئله باعث می‌شود که LoRA در مدت‌زمان کوتاهی تنظیم شود و کیفیت بالای خروجی‌ها را حفظ کند. متن‌های تولیدشده توسط مدل‌های فاین‌تیون شده با LoRA روان، منسجم، بدون خطاهای گرامی و کاملاً مرتبط با پرامپت و روودی هستند. این روش همچنین انعطاف‌پذیری بالایی دارد و می‌توان از آن برای بهینه‌سازی مدل‌های زبانی مختلف بدون نیاز به تغییر معماری اصلی مدل استفاده کرد. یکی دیگر از مزایای LoRA این است که وزن‌های اصلی مدل را حفظ می‌کند و تنها یک مجموعه‌ی کوچک از وزن‌های جدید به ساختار مدل اضافه می‌شود که این امر باعث

می شود مدل همچنان قابلیت استفاده از دانش قبلی خود را داشته باشد. این ویژگی LoRA را به بهترین گزینه برای تنظیم مدل های زبانی بزرگ تبدیل می کند. با این حال، یکی از چالش های LoRA این است که در برخی وظایف خاص، برای دستیابی به بهترین نتایج نیاز به انتخاب صحیح پارامترهای آن مانند مقدار rank و ماتریس های کمرتبه دارد. همچنین، پیاده سازی این روش نیازمند دانش تخصصی در حوزه ای بهینه سازی مدل های زبانی و نحوه ای اعمال تغییرات در لایه های خاص مدل است.

در مقایسه با LoRA، روش PEFT نیز عملکرد نسبتاً خوبی دارد اما همچنان از برخی جهات بهینه ترین روش محسوب نمی شود. یکی از مزایای PEFT این است که مصرف منابع آن کمتر از Fine-Tuning سنتی است، زیرا برخلاف روش های سنتی، نیازی به تنظیم تمام وزن های مدل ندارد و از تکنیک هایی مانند PEFT و Prefix-Tuning و Prompt-Tuning معمولاً خوب است و پاسخ های تولید شده دارای ساختار زبانی مناسبی هستند، اما در برخی موارد این روش دچار مشکلاتی در انتخاب واژگان یا ارتباط معنایی جمله ها می شود. یکی دیگر از مزایای PEFT این است که انعطاف پذیری نسبتاً بالایی دارد و می توان از آن برای تنظیم مدل های مختلف و بهینه سازی آن ها در وظایف گوناگون استفاده کرد. همچنین، این روش مانند LoRA، وزن های اصلی مدل را حفظ می کند و تنها تغییرات محدودی روی مدل اعمال می شود که این امر باعث می شود مدل بتواند هم از داده های جدید استفاده کند و هم اطلاعات قبلی خود را حفظ کند. یکی دیگر از مزایای PEFT این است که در برخی سناریوهای می توان بدون نیاز به فاین تیون کامل مدل، آن را برای وظایف خاص، مدل های فاین تیون شده با این روش ممکن است پاسخ هایی تولید کنند که ارتباط کاملی با پرامپت ورودی ندارند. همچنین برخی تکنیک های PEFT مانند Adapter-Tuning نیاز به تغییرات خاص در مدل دارند که این امر پیاده سازی آن را پیچیده تر می کند.

در مقایسه با LoRA و PEFT، روش Fine-Tuning دستی که در آن تمام لایه های مدل به جز لایه ای اول و آخر فریز شده بودند، ضعیف ترین عملکرد را داشت و معايب زیادی نسبت به دو روش دیگر داشت. یکی از محدود مزایای این روش این است که می توان مدل را برای داده های خاص بهینه کرد و در برخی موارد، تنظیم دستی مدل می تواند کمک کند تا مدل دانش جدیدی کسب کند. همچنین، در این روش کنترل کاملی روی فرآیند تنظیم مدل وجود دارد و می توان تمام تغییرات را به طور دقیق مدیریت کرد. اما این روش دارای مشکلات متعددی است که باعث شده عملکرد آن در مقایسه با LoRA و PEFT بسیار ضعیفتر باشد. یکی از بزرگ ترین معايب Fine-Tuning دستی، میزان بالای مصرف منابع محاسباتی است. این روش

نیاز به پردازش بسیار سنگینی دارد و مصرف حافظه‌ی GPU در آن بسیار بالا است، زیرا هر بار که مدل پردازش انجام می‌دهد، تمام وزن‌های مربوط به لایه‌های اول و آخر باید به روزرسانی شوند. علاوه بر این، این روش کندترین روش آموزش مدل محسوب می‌شود و در برخی موارد، مدت زمان موردنیاز برای فاین‌تیون مدل چندین برابر بیشتر از LoRA و PEFT است. یکی دیگر از مشکلات این روش، کیفیت پایین خروجی‌های مدل فاین‌تیون شده است. مدل‌هایی که با این روش تنظیم شده‌اند، عموماً دچار تکرارهای بی‌دلیل در پاسخ‌ها می‌شوند، درک نادرستی از پرامپت دارند، جملات بی‌معنی تولید می‌کنند و ترکیب زبان‌ها در آن‌ها به درستی رعایت نشده است. این مشکل به دلیل آن است که تنها دو لایه‌ی اول و آخر مدل تنظیم شده‌اند و سایر لایه‌ها فریز شده‌اند که این امر باعث شده مدل نتواند به درستی یادگیری کند و تنظیمات انجام‌شده روی وزن‌های جدید تأثیر کافی در عملکرد مدل نداشته باشد. یکی دیگر از معایب Fine-Tuning دستی، عدم انعطاف‌پذیری مناسب آن است. برخلاف LoRA و PEFT که انعطاف‌پذیری بالایی در تنظیم مدل‌های مختلف دارند، روش Fine-Tuning دستی نیازمند تغییرات گسترده در معماری مدل است و در بسیاری از موارد امکان بهینه‌سازی مدل برای وظایف جدید وجود ندارد. علاوه بر این، در این روش احتمال Overfitting افزایش پیدا می‌کند، زیرا مدل بیش از حد روی داده‌های جدید تنظیم می‌شود و ممکن است توانایی تعمیم‌دهی خود را از دست بدهد. یکی دیگر از مشکلات این روش این است که به دلیل فریز شدن بیشتر لایه‌های مدل، اطلاعات دریافتی از داده‌های آموزشی جدید به درستی در سراسر مدل گسترش پیدا نمی‌کند و مدل نمی‌تواند خود را برای وظایف جدید به درستی تنظیم کند.

در نهایت، مقایسه‌ی مزايا و معایب این روش‌ها نشان داد که LoRA بهترین روش برای Fine-Tuning مدل‌های زبانی بزرگ است، زیرا مصرف منابع آن کم است، سریع‌ترین زمان آموزش را دارد، کیفیت خروجی‌های آن بسیار بالا است و مدل می‌تواند خود را به طور دقیق برای داده‌های جدید تنظیم کند. روش PEFT نیز گزینه‌ی نسبتاً مناسبی است و عملکرد خوبی دارد اما همچنان در برخی موارد نیاز به تنظیمات پیچیده‌تر دارد و کیفیت خروجی‌های آن ممکن است به اندازه‌ی LoRA مطلوب نباشد. در مقابل، روش Fine-Tuning دستی که در آن تنها لایه‌ی اول و آخر مدل تنظیم شده‌اند، ناکارآمدترین روش محسوب می‌شود، زیرا مصرف حافظه‌ی زیادی دارد، زمان آموزش آن بسیار طولانی است، کیفیت خروجی‌های آن پایین است و مدل درک درستی از پرامپت ارائه نمی‌دهد. در مجموع، LoRA بهترین روش برای تنظیم مدل‌های زبانی است، PEFT در جایگاه دوم قرار دارد و روش Fine-Tuning دستی به دلیل مصرف بالای منابع و ضعف‌های جدی در عملکرد مدل، ناکارآمدترین روش محسوب می‌شود.

بر اساس تحلیلهای انجام‌شده، بهترین روش را با توجه به شرایط خاص تمرين معرفی و نتیجه‌گیری کنيد

با توجه به تمامی تحلیل‌های انجام‌شده در بخش‌های مختلف این تمرین، مشخص شد که روش LoRA بهترین گزینه برای Fine-Tuning مدل‌های زبانی در این سناریو است. این نتیجه‌گیری بر اساس معیارهای متعددی از جمله مصرف منابع، کیفیت خروجی‌ها، سرعت پردازش، میزان درک مدل از پرامپت، دقت پاسخ‌های تولیدی، انعطاف‌پذیری روش و هزینه‌ی محاسباتی انجام شده است. در مقایسه‌ی این روش با دو روش دیگر یعنی PEFT و LoRA نه تنها کمترین میزان منابع محاسباتی را مصرف می‌کند، بلکه بالاترین کیفیت پاسخ‌های خروجی را نیز ارائه می‌دهد.

یکی از اصلی‌ترین دلایلی که LoRA را به بهترین روش در این تمرین تبدیل کرده است، میزان مصرف بهینه‌ی منابع و سرعت بالای آموزش مدل است. برخلاف روش Fine-Tuning دستی که در آن تنها دو لایه‌ی مدل باز گذاشته شدند و نیاز به پردازش سنگین روی لایه‌های خاصی از مدل وجود داشت، LoRA با استفاده از ماتریس‌های کم‌رتبه، تغییرات خود را تنها روی بخش‌های موردنیاز اعمال می‌کند و در نتیجه، بدون نیاز به پردازش سنگین، مدل را با دقت بالایی تنظیم می‌کند. این مسئله باعث می‌شود که زمان پردازش بهشت‌کاهش پیدا کند، مصرف حافظه‌ی GPU بسیار پایین‌تر از سایر روش‌ها باشد و در عین حال، خروجی‌های مدل از نظر دقت و انسجام بالاترین کیفیت را داشته باشند. در مقابل، روش Fine-Tuning دستی بیشترین میزان مصرف منابع را داشت، زمان آموزش آن طولانی بود و مدل نتوانست پاسخ‌های مناسبی تولید کند. همچنین، روش PEFT در این مقایسه عملکرد نسبتاً خوبی داشت اما همچنان در برخی موارد، میزان پردازش آن بیشتر از LoRA بود و کیفیت خروجی‌های آن به اندازه‌ی LoRA مطلوب نبود.

از نظر کیفیت خروجی‌ها و میزان درک مدل از پرامپت‌های ورودی، LoRA عملکرد بهتری نسبت به سایر روش‌ها داشت. این روش توانست پاسخ‌هایی تولید کند که از نظر زبان‌شناسی روان، بدون تکرار بی‌دلیل، دارای ساختار صحیح، مرتبط با پرامپت و بدون ترکیب نامناسب زبان‌ها بودند. در مقایسه، روش Fine-Tuning دستی عملکرد بسیار ضعیفی داشت و مدل دچار مشکلات متعددی از جمله درک نادرست از پرامپت، تکرار بی‌دلیل کلمات، ترکیب نامناسب فارسی و انگلیسی و ارائه‌ی پاسخ‌های غیرقابل فهم شد. این مسئله نشان می‌دهد که محدود کردن فرآیند یادگیری به دو لایه‌ی اول و آخر مدل باعث کاهش توانایی یادگیری مدل شده است و مدل نتوانسته خود را برای داده‌های جدید به درستی تنظیم کند. از سوی دیگر، روش PEFT عملکرد متوسطی داشت و اگرچه پاسخ‌های آن در بسیاری از موارد مناسب بودند، اما همچنان در برخی نمونه‌ها مشکلات گرامری و معنایی مشاهده شد. این مسئله نشان داد که LoRA

می‌تواند بالاترین کیفیت خروجی را ارائه دهد و پاسخ‌های آن درک بسیار بهتری از پرامپت ورودی داشته باشند.

یکی دیگر از دلایل انتخاب LoRA به عنوان بهترین روش در این تمرین، انعطاف‌پذیری بالای آن در تنظیم مدل‌های زبانی مختلف است. این روش برخلاف Fine-Tuning دستی که نیاز به تنظیم مستقیم وزن‌های مدل دارد و ممکن است باعث بروز مشکلاتی مانند Overfitting شود، به گونه‌ای طراحی شده که می‌توان آن را روی مدل‌های مختلف پیاده‌سازی کرد و بدون تغییر در معماری اصلی مدل، مدل را برای وظایف مختلف تنظیم کرد. این انعطاف‌پذیری به ویژه در شرایطی که نیاز به بهینه‌سازی مدل برای وظایف خاص بدون تغییرات اساسی در معماری آن وجود دارد، یک مزیت بسیار مهم محسوب می‌شود. در مقابل، روش Fine-Tuning دستی به دلیل محدود بودن تغییرات تنها به دو لایه‌ی خاص، نتوانست تطبیق لازم را با وظایف جدید پیدا کند و باعث شد که مدل یادگیری ناقصی داشته باشد. همچنین، روش PEFT نیز انعطاف‌پذیری بالایی دارد اما در برخی موارد، به تنظیمات پیچیده‌تر نیاز دارد و عملکرد آن به اندازه‌ی بهینه LoRA نیست.

از نظر هزینه‌ی پردازشی و میزان مصرف منابع، LoRA باز هم برترین روش محسوب می‌شود. این روش کمترین میزان حافظه‌ی GPU را مصرف کرد، سریع‌ترین زمان پردازش را داشت و بدون نیاز به تنظیم کل پارامترهای مدل، توانست خروجی‌های مناسبی تولید کند. در مقابل، روش Fine-Tuning دستی نه تنها از نظر هزینه‌ی پردازشی ناکارآمد بود، بلکه نیاز به زمان طولانی برای پردازش داشت و در نهایت هم نتوانست کیفیت مناسبی در پاسخ‌ها ارائه دهد. روش PEFT نیز میزان مصرف منابع کمتری نسبت به Fine-Tuning داشت اما همچنان در مقایسه با LoRA نیازمند پردازش بیشتری بود و در برخی موارد، به دلیل روش‌های مورد استفاده در Prompt-Tuning، میزان مصرف حافظه‌ی آن بالاتر از LoRA بود. بنابراین، LoRA از نظر هزینه‌ی پردازشی و هم از نظر کیفیت خروجی، بهینه‌ترین روش است و در این تمرین بهترین عملکرد را ارائه داد.

در نهایت، با در نظر گرفتن تمام این تحلیل‌ها، مشخص شد که LoRA بهترین روش برای تنظیم مدل‌های زبانی در این سناریو است، زیرا توانست با کمترین میزان مصرف منابع، بالاترین کیفیت خروجی را ارائه دهد، فرآیند آموزش را در کوتاه‌ترین زمان ممکن انجام دهد، بدون تغییر در معماری اصلی مدل، آن را برای وظایف جدید تنظیم کند و انعطاف‌پذیری بالایی برای استفاده در مدل‌های مختلف داشته باشد. روش PEFT نیز گزینه‌ی نسبتاً مناسبی بود اما همچنان در برخی موارد میزان مصرف منابع و کیفیت خروجی‌های آن نسبت به LoRA پایین‌تر بود. در مقابل، روش Fine-Tuning دستی که در آن فقط دو لایه‌ی اول و آخر مدل آزاد شده بودند، ضعیفترین روش بود، زیرا مصرف منابع آن بالا بود، مدل به درستی

آموزش ندید، خروجی‌های نامناسب و غیرمنسجمی تولید کرد و نتوانست عملکرد خوبی از خود نشان دهد. بنابراین، LoRA بهترین روش، PEFT گزینه‌ای قابل قبول اما کمتر بهینه و Fine-Tuning دستی ناکارآمدترین روش در این تمرین محسوب می‌شود.

پرسش ۲ - تولید کپشن برای تصاویر

۱-۱. آماده سازی دیتاست



a dog runs across the grass .
A golden dog is running on the grass .
A tan dog leaping in a grassy field .
Brown dog leaps through field .
The dog is running quickly through the meadow .



A man is sitting on a black and brown dog .
A man stands over a brown and black dog in the grass .
A man wearing a black and white shirt jumping on a dog 's back .
An attack dog is working with his trainer .
A running dog and a standing man on a dry field of grass .



A boy in a black cap break dancing
A young adult dances on his hand .
Guy break dancing on brick floor in front of students ' union .
One man watches while another breakdances , balancing himself for a moment on one hand .
The male teenager is break dancing on the red and black tiled floor .



Two black dogs are facing each other showing their teeth .
two brown and black dogs growling at each other .
Two brown dogs barking at each other .
Two dogs bare their teeth and bark .
Two dogs on the grass , barking at each other .



Two girls are walking down a dirt road in a park .
Two women walking down a dirt road with a man far behind .
Two women walking on an outdoor trail .
Two young women are walking along a rural path bordered by trees .
Two young women walk along a dirt road .

شکل ۴ نمونه هایی از تصاویر خوانده شده به همراه کپشن ها

برای سازگاری با ورودی مدل CNN، ابعاد تصاویر را به ابعاد ثابت و مناسب تغییر دادیم. به طور معمول، ابعاد ۲۲۴*۲۲۴ پیکسل برای مدل‌های CNN مانند VGG16 و ResNet مناسب است. این کار به ما کمک می‌کند تا از پردازش‌های اضافی و زمان بر جلوگیری کنیم و همچنین از مشکلاتی مانند تغییر مقیاس و نسبت ابعاد جلوگیری شود.

نرمالسازی یکی از مراحل کلیدی در پیش‌پردازش تصاویر است. در این مرحله، مقادیر پیکسل‌های تصاویر را به محدوده‌ای خاص نرمالسازی کردیم. ما از روش نرمالسازی با استفاده از میانگین و انحراف معیار استاندارد استفاده کردیم.



شکل ۵ نمونه تصاویر پیش‌پردازش شده با کپشن پیش‌پردازش شده

پیش‌پردازش متون

برای شروع، متون توصیفی را به حروف کوچک تبدیل کردم تا حساسیت به حروف بزرگ حذف شود. همچنین، علائم نگارشی، نمادهای خاص و اعداد غیرضروری را حذف کردم. این کار به ما کمک می‌کند تا فقط بر روی کلمات اصلی تمرکز کنیم و از نویزهای اضافی جلوگیری کنیم. توکنیزاسیون

پس از پیش‌پردازش، مرحله توکنیزاسیون آغاز شد. در این مرحله، هر کلمه به یک شناسه عددی منحصر به فرد تبدیل شد. برای این کار، یک دیکشنری ایجاد کردم که در آن هر کلمه به یک عدد نگاشت شده است. همچنین، special token های مورد نیاز مانند <sos>, <eos> و <junk> نیز به این دیکشنری اضافه کردم.

پس از توکنیزاسیون، برای اطمینان از اینکه تمام توکن‌ها دارای طول یکسانی هستند، از پدینگ استفاده کردم. طول حداکثر توکن‌ها را به ۱۵ محدود کردم و اگر تعداد توکن‌ها کمتر از این مقدار بود، با استفاده از <pad> آن‌ها را پر کردم

توضیح دهید که چرا باید کپشنها را به طول ثابت تبدیل کنیم؟

۱. سازگاری با ورودی مدل

مدل‌های یادگیری عمیق، به ویژه شبکه‌های عصبی کانولوشنی (CNN) و شبکه‌های عصبی بازگشتی (RNN)، نیاز دارند که ورودی‌های آن‌ها دارای ابعاد ثابتی باشند. در تسك ایمیج کپشنینگ، کپشن‌ها به عنوان ورودی برای مدل‌های تولید متن عمل می‌کنند. اگر طول کپشن‌ها متفاوت باشد، مدل نمی‌تواند به درستی آن‌ها را پردازش کند و این موضوع می‌تواند منجر به بروز خطاهای اجرایی و عدم کارایی مدل شود. بنابراین، تبدیل کپشن‌ها به طول ثابت به مدل این امکان را می‌دهد که ورودی‌ها را به صورت یکنواخت دریافت کند و از این رو، به بهبود عملکرد آن کمک می‌کند.

۲. بهبود کارایی محاسباتی

تبدیل کپشن‌ها به طول ثابت می‌تواند به بهینه‌سازی محاسبات کمک کند. با داشتن ورودی‌های یکسان، می‌توان از تکنیک‌های مختلفی برای بهبود سرعت محاسبات و کاهش زمان آموزش استفاده کرد. به عنوان مثال، در هنگام آموزش مدل، داده‌ها معمولاً به صورت دسته‌ای (batch) پردازش می‌شوند. اگر طول کپشن‌ها متفاوت باشد، پردازش دسته‌ای می‌تواند دشوار و زمان بر باشد. با استفاده از پدینگ و تبدیل کپشن‌ها به طول ثابت، می‌توان به راحتی داده‌ها را در دسته‌های یکنواخت تقسیم کرد و این امر به بهبود کارایی و سرعت آموزش کمک می‌کند.

۳. مدیریت داده‌های گمشده

در بسیاری از موارد، کپشن‌ها ممکن است کوتاه‌تر از طول مورد نظر باشند. در چنین شرایطی، استفاده از پدینگ به ما این امکان را می‌دهد که این کپشن‌ها را به طول ثابت تبدیل کنیم. این کار به مدل کمک می‌کند تا با داده‌های گمشده یا ناقص بهتر برخورد کند و از تأثیر منفی آن‌ها بر روی یادگیری جلوگیری کند. به عبارت دیگر، با استفاده از پدینگ، می‌توانیم از اطلاعات موجود در کپشن‌ها به بهترین نحو استفاده کنیم و از بروز مشکلات ناشی از داده‌های ناقص جلوگیری نماییم.

۴. یکپارچگی داده‌ها

تبدیل کپشن‌ها به طول ثابت به یکپارچگی داده‌ها کمک می‌کند. این یکپارچگی به مدل اجازه می‌دهد تا الگوهای موجود در داده‌ها را بهتر شناسایی کند و به یادگیری عمیق‌تر و مؤثرتری منجر شود. به عنوان مثال، اگر یک مدل با داده‌های یکپارچه و استاندارد آموزش ببیند، احتمالاً می‌تواند الگوهای پیچیده‌تری را در داده‌ها شناسایی کند و در نتیجه، عملکرد بهتری در تولید کپشن‌های متناسب با تصاویر داشته باشد.

۵. کاهش پیچیدگی در پردازش زبان طبیعی

در تسک ایمیج کپشنینگ، تولید متن به عنوان یک پروسه پیچیده در نظر گرفته می‌شود. با تبدیل کپشن‌ها به طول ثابت، می‌توان پیچیدگی‌های پردازش زبان طبیعی را کاهش داد. این امر به مدل کمک می‌کند تا به راحتی با ساختارهای زبانی مختلف و الگوهای گرامری آشنا شود و در نتیجه، توانایی تولید متن‌های روان و طبیعی را افزایش دهد.

۶. تسهیل در یادگیری انتقالی

مدل‌های یادگیری عمیق معمولاً از یادگیری انتقالی Transfer Learning بهره می‌برند، به این معنا که از مدل‌های پیش‌آموزش‌دیده استفاده می‌کنند و آن‌ها را برای تسک‌های خاص خود تنظیم می‌کنند. با تبدیل کپشن‌ها به طول ثابت، می‌توان از مدل‌های پیش‌آموزش‌دیده به طور مؤثرتری استفاده کرد، زیرا این مدل‌ها معمولاً با داده‌هایی با طول ثابت آموزش دیده‌اند. این امر می‌تواند به تسريع در فرایند آموزش و بهبود عملکرد نهایی مدل کمک کند.

تقسیم دادگان

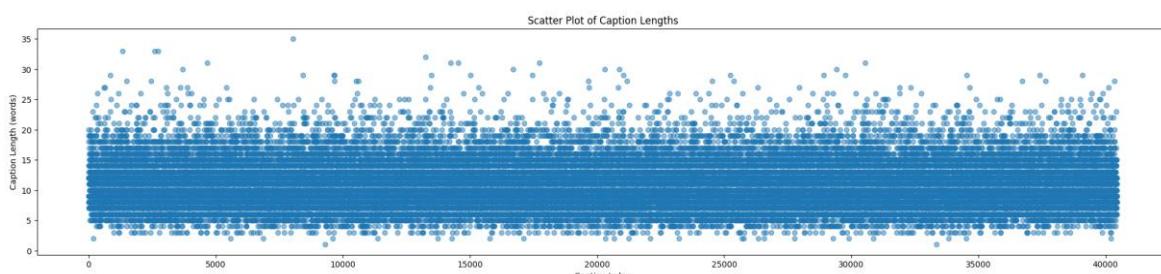
در این بخش از پژوهه، به تقسیم داده‌ها به مجموعه‌های آموزشی، اعتبارسنجی و آزمایش پرداختم. این مرحله یکی از مراحل کلیدی در فرآیند یادگیری عمیق است، زیرا به ما این امکان را می‌دهد که مدل را بر روی یک مجموعه داده آموزش دهیم و سپس عملکرد آن را بر روی مجموعه‌های دیگر ارزیابی کنیم. این

کار به ما کمک می کند تا از بروز مشکلاتی مانند overfitting جلوگیری کنیم و اطمینان حاصل کنیم که مدل به درستی تعمیم‌پذیری دارد.

برای تقسیم داده‌ها، از تابع `train_test_split` از کتابخانه `sklearn.model_selection` استفاده کردم. این تابع به ما این امکان را می‌دهد که به راحتی داده‌ها را به چندین مجموعه تقسیم کنیم. در ابتدا، تمامی تصاویر پیش‌پردازش شده را به دو دسته تقسیم کردم: مجموعه آموزشی و مجموعه موقت. در اینجا، ۸۰ درصد از داده‌ها به مجموعه آموزشی و ۲۰ درصد به مجموعه موقت اختصاص داده شد. پس از آن، مجموعه موقت را به دو قسمت تقسیم کردم: مجموعه اعتبارسنجی و مجموعه آزمایش. به این ترتیب، ۱۰ درصد از کل داده‌ها به هر یک از این دو مجموعه اختصاص یافت. این تقسیم‌بندی به ما این امکان را می‌دهد که مدل را بر روی مجموعه آموزشی آموزش دهیم و سپس با استفاده از مجموعه اعتبارسنجی، پارامترهای مدل را بهینه کنیم. در نهایت، از مجموعه آزمایش برای ارزیابی نهایی عملکرد مدل استفاده خواهیم کرد.

پس از تقسیم داده‌ها، جداول‌های DataFrame‌هایی برای هر یک از مجموعه‌ها ایجاد کردم. این کار به ما کمک می‌کند تا به راحتی به داده‌های هر مجموعه دسترسی داشته باشیم و آن‌ها را برای مراحل بعدی پردازش کنیم. به عنوان مثال، DataFrame مجموعه آموزشی شامل تمام سطرهایی است که تصاویر آن‌ها در مجموعه آموزشی قرار دارند. در نهایت، تعداد نمونه‌ها در هر یک از مجموعه‌ها را چاپ کردم تا از صحت تقسیم داده‌ها اطمینان حاصل شود. نتایج نشان می‌دهد که ۶۴۷۲ نمونه به مجموعه آموزشی، ۸۰۹ نمونه به مجموعه اعتبارسنجی و ۸۱۰ نمونه به مجموعه آزمایش اختصاص داده شده است. این تقسیم‌بندی به ما این اطمینان را می‌دهد که مدل به طور مؤثری آموزش دیده و توانایی تعمیم‌پذیری به داده‌های جدید را دارد.

نمودار پراکندگی (Scatter Plot) طول کپشن‌ها را رسم کنید تا تنوع طول توصیفها در دیتابست مشخص شود

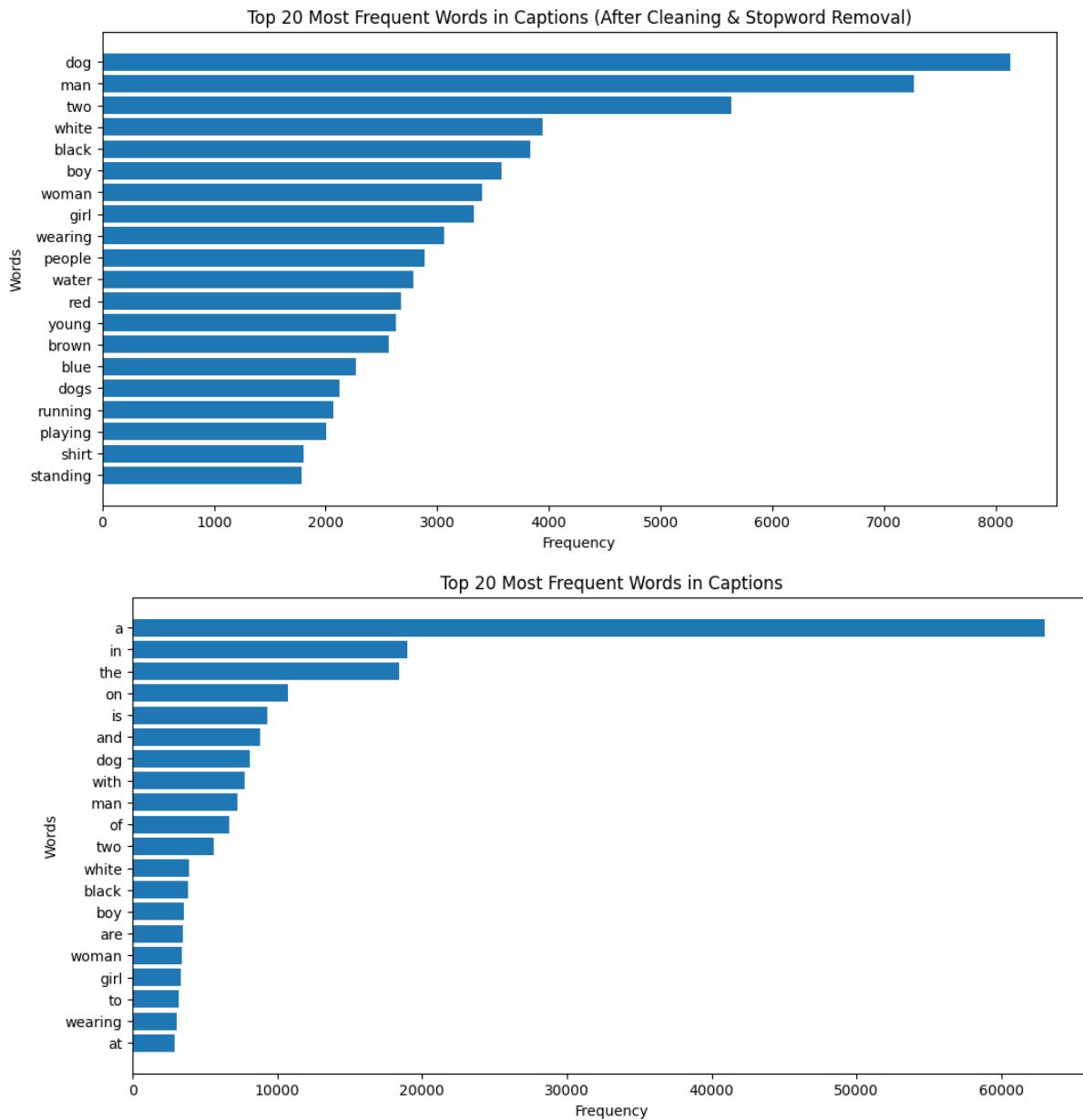


شکل ۶ پراکندگی طول کپشن‌ها

پس از محاسبه طول کپشن‌ها، یک نمودار پراکندگی (scatter plot) برای نمایش این داده‌ها ترسیم کردم. برای این کار از کتابخانه `matplotlib` استفاده کردم و اندازه شکل نمودار را به ۲۵ در ۵ تنظیم کردم

تا بتوانم به وضوح توزیع طول کپشن‌ها را مشاهده کنم. در محور افقی، ایندکس کپشن‌ها و در محور عمودی، طول آن‌ها به صورت تعداد کلمات نمایش داده شد. با استفاده از تابع (`scatter`)، نقاط مربوط به هر کپشن در نمودار نمایش داده شدند و شفافیت نقاط را با استفاده از پارامتر `alpha` به ۰,۵ تنظیم کردم تا نقاط تداخل نکنند و بهتر قابل مشاهده باشند. در نهایت، عنوان و برچسب‌های محورهای نمودار را اضافه کردم و نمودار را با استفاده از تابع (`show`) نمایش دادم.

پس از نمایش نمودار، اطلاعاتی را در مورد آن ثبت کردم و مدت زمان اجرای این بخش از کد را نیز نمایش دادم. نتایج نشان می‌دهد که بیشتر داده‌های ما در بازه طولی ۷ تا ۱۸ کلمه قرار دارند. این نکته حائز اهمیت است زیرا نشان می‌دهد که اکثر کپشن‌ها دارای طولی مناسب و قابل قبول هستند. همچنین، ماسکیمم طول کپشن‌ها ۳۵ کلمه و مینیمم آن‌ها ۳ کلمه است. این اطلاعات می‌تواند به ما در انتخاب پارامترهای مناسب برای مدل و همچنین در طراحی و بهینه‌سازی الگوریتم‌های پردازش زبان طبیعی کمک کند. به همین دلیل، هنگام استفاده از روش‌های جستجو و انتخاب، از این بازه‌های طولی استفاده خواهیم کرد تا بهترین نتایج را به دست آوریم.



شکل ۷ تعداد تکرار کلمات در کپشن ها

در این پژوهه ما به تحلیل و پردازش کپشن‌های مربوط به مجموعه‌ای از تصاویر پرداختیم. هدف اصلی این بخش از پژوهه، بررسی و استخراج پرترکارترین کلمات مورد استفاده در کپشن‌ها و همچنین حذف کلمات زائد (Stopwords) برای تمرکز بر کلمات کلیدی و مهم است. در واقع، پردازش زبان طبیعی (NLP) یکی از مهم‌ترین مراحل در تحلیل متون متنی مرتبط با تصاویر است، زیرا اطلاعات استخراج شده از کپشن‌ها می‌تواند تأثیر زیادی در بهبود مدل‌های هوش مصنوعی که برای توصیف تصاویر استفاده می‌شوند، داشته باشد. به همین دلیل، ما تصمیم گرفتیم که این تحلیل را در دو مرحله انجام دهیم: ابتدا بررسی تمام کلمات بدون هیچ‌گونه پردازشی، و سپس حذف کلمات زائد که معمولاً ارزش معنایی خاصی ندارند و فقط نقش گرامی دارند.

در مرحله اول، ما تمام واژه‌های موجود در کپشن‌های مجموعه داده را جمع‌آوری کرده و سپس تعداد تکرار هر واژه را محاسبه کردیم تا پرتکرارترین کلمات را مشخص کنیم. نتیجه این مرحله را در اولین نمودار می‌توان مشاهده کرد. نکته‌ای که در این نمودار به چشم می‌خورد این است که بسیاری از کلمات پرتکرار شامل واژه‌هایی مانند "a" ، "the" ، "in" ، "on" ، "is" ، "and" و "of" هستند. این واژه‌ها عموماً در هر متن انگلیسی پرتکرار هستند زیرا به ساختار جمله کمک می‌کنند، اما در واقع محتوای خاصی به ما ارائه نمی‌دهند و اطلاعات مفیدی درباره موضوع کپشن‌ها در اختیارمان نمی‌گذارند. در حقیقت، این کلمات بیشتر نقش اتصال‌دهنده دارند تا این که نشان‌دهنده‌ی یک مفهوم خاص باشند. بنابراین، مرحله دوم پردازش داده‌ها که شامل حذف این کلمات بود، اهمیت زیادی پیدا کرد.

در مرحله دوم، ما از یک مجموعه از کلمات توقف (Stopwords) که در زبان انگلیسی جزو واژه‌های رایج ولی بی‌اهمیت محسوب می‌شوند، استفاده کردیم و این کلمات را از لیست واژه‌های پرتکرار حذف کردیم. نتیجه این فرآیند را می‌توان در نمودار دوم مشاهده کرد. این نمودار نشان می‌دهد که بعد از حذف این کلمات زائد، پرتکرارترین کلمات باقی‌مانده شامل مواردی مانند "dog" ، "man" ، "woman" ، "boy" و "girl" هستند. نکته جالبی که در این نتایج به چشم می‌خورد این است که برخلاف مرحله قبل، اکنون ما به جای کلمات گرامری، با کلمات مفهومی و پرمونا روبرو هستیم. این نشان می‌دهد که پردازش انجام‌شده توانسته است داده‌های ما را تمیزتر و کاربردی‌تر کند.

برای مثال، مشاهده می‌کنیم که کلمه "dog" یکی از پرتکرارترین کلمات است. این نشان می‌دهد که در بسیاری از کپشن‌های موجود در این مجموعه داده، سگ به عنوان یک سوژه‌ی اصلی در تصویر حضور دارد. از سوی دیگر، کلماتی مانند "woman" و "man" نیز بسیار پرتکرار هستند که به این معناست که بسیاری از تصاویر شامل انسان‌ها می‌شوند. این اطلاعات می‌تواند در بهبود مدل‌های توصیف تصویر کمک زیادی کند، چرا که به ما می‌گوید که در این مجموعه داده، چه چیزهایی بیشتر دیده می‌شود و مدل باید روی چه مفاهیمی بیشتر تمرکز کند. علاوه بر این، ما مشاهده می‌کنیم که برخی از رنگ‌ها مانند "white" ، "black" ، "blue" ، "red" و "brown" نیز در لیست پرتکرارترین کلمات قرار دارند که نشان می‌دهد کپشن‌های ما به رنگ اشیاء یا لباس‌های موجود در تصویر اشاره دارند. این موضوع نیز از این جهت اهمیت دارد که می‌تواند در مدلسازی ویژگی‌های بصری تصاویر مفید باشد.

همچنین، افعالی مانند "running" و "playing" در بین پر تکرارترین کلمات قرار دارند که نشان می‌دهد کپشن‌های مربوط به این مجموعه داده نه تنها به اشیاء و موجودات زنده اشاره دارند، بلکه به فعالیت‌هایی که در تصویر در حال انجام هستند نیز می‌پردازند. این موضوع می‌تواند برای مدل‌های یادگیری عمیق که هدف‌شان تشخیص فعالیت‌ها در تصاویر است، بسیار کاربردی باشد.

در نهایت، این تحلیل نشان داد که اگرچه در مرحله‌ی اول بسیاری از کلمات پر تکرار مربوط به ساختار زبانی و گرامری بودند و ارزش معنایی خاصی نداشتند، اما پس از حذف کلمات زائد توانستیم کلمات کلیدی واقعی را شناسایی کنیم که به درک بهتر محتوای تصاویر کمک می‌کنند. نتایج این تحلیل همچنین نشان می‌دهند که در این مجموعه داده، بسیاری از تصاویر شامل انسان‌ها، حیوانات (مخصوصاً سگ‌ها)، لباس‌ها و رنگ‌ها هستند و بسیاری از کپشن‌ها نیز به فعالیت‌های فیزیکی مانند دویدن و بازی کردن اشاره دارند.

با این تحلیل، ما اکنون درکی دقیق‌تر از داده‌های خود داریم که می‌تواند در مراحل بعدی پروژه، مانند مدل‌سازی و بهبود الگوریتم‌های کپشن‌گذاری تصاویر، مفید باشد. این اطلاعات می‌توانند به بهینه‌سازی مدل‌های یادگیری ماشین و بهبود کیفیت کپشن‌های تولید شده کمک کنند، زیرا اکنون می‌دانیم که چه کلمات کلیدی و چه مفاهیمی بیشتر در داده‌های ما دیده می‌شوند و در نتیجه می‌توانیم مدل را به گونه‌ای تنظیم کنیم که بتواند این الگوها را بهتر یاد بگیرد.

۲.۳ پیاده سازی CNN-RNN

در این بخش از پروژه، ما مدل CNN-RNN را برای انجام توضیح تصویر (Image Captioning) پیاده‌سازی کردیم. این مدل از ترکیب دو نوع شبکه عصبی استفاده می‌کند: یک شبکه عصبی پیچشی (RNN) برای استخراج ویژگی‌های تصویری و یک شبکه عصبی بازگشتی (CNN) که مسئول تولید متن توصیفی مرتبط با تصویر است. به طور دقیق‌تر، ما از مدل EfficientNet-B0 به عنوان رمزگذار (Encoder) و از LSTM به عنوان رمزگشا (Decoder) استفاده کردیم. مدل EfficientNet-B0 یکی از مدل‌های از پیش‌آموزش داده شده است که دقت بالایی دارد و بهینه‌شده برای کارهای بینایی ماشین است. در ابتدا، لایه‌ی Fully Connected نهایی این مدل را حذف کردیم تا فقط ویژگی‌های بصری استخراج شوند. این ویژگی‌ها به عنوان ورودی رمزگشا استفاده می‌شوند تا کلمات توصیفی تصویر را تولید کند. این بردار ویژگی از طریق یک لایه خطی (Fully Connected Layer) به ابعاد مناسب کاهش پیدا کرد و سپس به LSTM داده شد. دلیل استفاده از LSTM این است که این مدل در پردازش داده‌های ترتیبی عملکرد بهتری نسبت به RNN معمولی دارد و مشکل از بین رفتن گرادیان را برطرف می‌کند.

در بخش رمزگشا (Decoder)، ابتدا یک لایه Embedding برای تبدیل کلمات به بردارهای عددی استفاده کردیم. دلیل این کار این است که نمایش کلمات در قالب بردارهای تعییه شده (Word Embedding) باعث می شود که مدل بتواند روابط معنایی بین کلمات را بهتر درک کند و به جای اینکه هر کلمه را به صورت یک بردار One-Hot نمایش دهد که بسیار پراکنده و ناکارآمد است، می توان آن را در یک فضای برداری متراکم و معنادار نمایش داد. سپس این بردارهای عددی همراه با ویژگی های تصویری وارد LSTM شدند تا مدل بتواند توصیف مربوط به تصویر را مرحله به مرحله تولید کند. بعد از خروجی گرفتن از LSTM، یک لایه‌ی خطی همراه با Softmax استفاده شد تا مدل بتواند احتمال وقوع هر کلمه را پیش‌بینی کند و کلمه‌ای که بیشترین احتمال را دارد، به عنوان خروجی انتخاب کند.

در گام بعد، برای اتصال رمزگذار و رمزگشا، ما یک کلاس ImageCaptioningModel پیاده‌سازی کردیم که شامل هر دو بخش CNN و LSTM بود. در این معماری، تصویر ابتدا از طریق رمزگذار پردازش شده و ویژگی های آن استخراج می شود. سپس، این ویژگی ها به رمزگشا داده می شوند تا متن توصیفی تولید شود. یکی از چالش های این مرحله، این بود که ویژگی های تصویری دارای ابعاد متفاوتی از بردارهای کلمات بودند. برای حل این مسئله، ویژگی های تصویری را از طریق یک لایه‌ی خطی به ابعاد مورد نیاز کاهش دادیم تا با بردارهای کلمات سازگار شود. در نهایت، مدل را به صورت End-to-End طراحی کردیم تا به راحتی قابل آموزش باشد.

برای آموزش مدل، ازتابع هزینه CrossEntropyLoss استفاده کردیم و مقدار Label Smoothing را برابر با $1 - \epsilon$ قرار دادیم تا مدل بیش از حد روی نمونه های خاص متمرکز نشود و تعمیم پذیری بهتری داشته باشد. همچنین، هنگام محاسبه هزینه، کلمات Padding را نادیده گرفتیم تا مدل تنها بر روی پیش‌بینی صحیح کلمات تمرکز کند. یکی از نکات مهم در آموزش این بود که همه‌ی لایه‌های CNN را فریز کردیم (به جز چند لایه آخر) تا ویژگی های تصویری از مدل از پیش‌آموزش داده شده گرفته شود و فقط لایه‌های بالایی برای این کار بهینه شوند. برای بهینه‌سازی، از AdamW استفاده کردیم که در مقایسه با عملکرد بهتری روی داده های تصویری دارد. همچنین، از ReduceLROnPlateau برای کاهش نرخ یادگیری در صورت عدم بهبود عملکرد استفاده کردیم تا مدل دچار نوسانات بی مورد نشود.

یکی از مهم‌ترین چالش‌ها در آموزش این مدل، جلوگیری از بیش‌برازش (Overfitting) بود. برای این کار، چندین تکنیک را به کار گرفتیم. نخست، در قسمت Data Augmentation، روی تصاویر ورودی تغییراتی مانند افقی‌سازی تصادفی (Random Horizontal Flip) و برش تصادفی (Random Resized Crop) اعمال کردیم تا مدل بهتر تعمیم یابد. دوم، از Label Smoothing در CrossEntropy استفاده کردیم تا مدل بیش از حد روی نمونه های خاص وابسته نشود. سوم، از Checkpoints برای ذخیره مدل در

فوائل مشخص و ادامه‌ی آموزش از همان نقطه استفاده کردیم تا در صورت نیاز بتوانیم مدل را از جایی که متوقف شده بود ادامه دهیم.

در بخش ارزیابی مدل، ما روند کاهش خطای مدل را در طول هر epoch بررسی کردیم تا ببینیم که آیا مدل در حال یادگیری است یا خیر. نمودار خطای داده‌های آموزش و ارزیابی را رسم کردیم تا مشخص شود که آیا مدل در حال بهبود است یا دچار بیش‌برازش شده است. همچنین، در پایان هر epoch، یک نمونه تصویر را گرفته و کپشن تولیدشده توسط مدل را نمایش دادیم تا ببینیم که عملکرد آن چگونه است. برای نمایش خروجی‌ها، از روش Beam Search استفاده کردیم که به جای انتخاب تنها محتمل‌ترین کلمه در هر مرحله، چندین کلمه را در نظر می‌گیرد و بهترین ترکیب از چندین کلمه را پیدا می‌کند. این کار باعث می‌شود که کپشن‌های تولیدی روان‌تر و معنادارتر باشند.

در نهایت، ما ۵ تصویر را همراه با کپشن‌های تولیدشده توسط مدل نمایش دادیم و عملکرد مدل را تحلیل کردیم. در برخی از موارد، مدل عملکرد خوبی داشت و توانست به درستی اشیاء و فعالیت‌ها را شناسایی کند. اما در برخی موارد، مدل اشیاء را به درستی تشخیص نمی‌داد و یا روابط بین آن‌ها را درست متوجه نمی‌شد. برای مثال، در برخی از کپشن‌های تولیدشده، مدل به اشتباہ یک سگ را به جای گربه تشخیص داده بود یا برخی فعالیت‌ها را به درستی بیان نکرده بود. این نشان داد که هنوز جا برای بهبود وجود دارد و برای دقت بالاتر می‌توان از مکانیزم Attention استفاده کرد که به مدل اجازه می‌دهد روی بخش‌های مهم تصویر تمرکز کند.

از یک لایه embedding برای کلمات استفاده کنید. دلیل این کار را توضیح دهید. چرا استفاده از تعبیه کلمات (Word Embedding) به جای بردارهای One-hot مناسب‌تر است؟

در پیاده‌سازی مدل توضیح تصویر از یک لایه تعبیه کلمات یا همان Word Embedding استفاده کردیم که این بخش یکی از اجزای بسیار مهم در شبکه عصبی ما محسوب می‌شود. این لایه کمک می‌کند که کلمات موجود در کپشن‌ها از حالت گسسته و مجزا که در روش One-Hot Encoding اتفاق می‌افتد، به نمایش برداری متراکم و معنادار تبدیل شوند. در واقع، در روش One-Hot Encoding هر کلمه در قالب یک بردار دودویی نمایش داده می‌شود که دارای طولی برابر با تعداد کل کلمات در مجموعه داده است و هر کلمه فقط یک مقدار یک در موقعیت مشخصی از بردار دارد و بقیه مقادیر صفر هستند. برای مثال، اگر واژگان ما شامل پنج کلمه باشد مانند cat، dog، apple، tree و banana، در روش One-Hot این کلمات به ترتیب در بردارهای مختلف به این شکل نمایش داده می‌شوند که cat معادل بردار $[1, 0, 0, 0, 0]$ است، dog معادل بردار $[0, 1, 0, 0, 0]$ ، apple معادل بردار $[0, 0, 1, 0, 0]$ ، tree معادل بردار $[0, 0, 0, 1, 0]$ و banana معادل بردار $[0, 0, 0, 0, 1]$ است. این نمایش اگرچه می‌تواند کلمات را به صورت

منحصر به فرد نمایش دهد اما هیچ اطلاعاتی درباره شباهت‌های معنایی میان آن‌ها ارائه نمی‌کند. در واقع مدل نمی‌تواند بفهمد که cat و dog از نظر مفهومی شباهت بیشتری به یکدیگر دارند تا cat و banana و این مسئله باعث می‌شود که مدل نتواند روابط معنایی بین کلمات را درک کند. علاوه بر این، زمانی که واژگان ما بسیار بزرگ شوند، بردارهای One-Hot بسیار طولانی خواهند شد و این موضوع از نظر مصرف حافظه و پردازش بسیار ناکارآمد است. برای مثال اگر ما پنجاه هزار کلمه در مجموعه داده خود داشته باشیم، هر کلمه باید در برداری با همین تعداد بعد نمایش داده شود که محاسبات را بسیار سنگین می‌کند و نیاز به منابع بسیار زیادی دارد.

در مقابل، روش Word Embedding این مشکل را برطرف می‌کند. در این روش، هر کلمه به یک بردار عددی فشرده و معنادار تبدیل می‌شود که معمولاً در فضایی با تعداد ابعاد کمتر مانند ۲۵۶ یا ۳۰۰ بعد نمایش داده می‌شود. این نمایش باعث می‌شود که کلماتی که از نظر معنایی به یکدیگر نزدیک هستند، در فضای برداری نیز به یکدیگر نزدیک باشند و در مقابل، کلمات نامرتبط از هم فاصله داشته باشند. به عنوان مثال، در یک مدل Word Embedding ممکن است بردارهای کلمات به این شکل باشند که cat معادل بردار $[0, 0, 12, 0, 0, 55, 0, 76, \dots]$ باشد، dog معادل بردار $[0, 1.1, 0, 0, 74, 0, 50, \dots]$ باشد، apple معادل بردار $[1, 23, 0, 55, 0, 0, 89, 0, 67.0, \dots]$ باشد، banana معادل بردار $[0, 0, 50, 0, 0, 65, 0, 0, 85, \dots]$ باشد و tree معادل بردار $[0, 45.0, 0, 89, 1, 23, 0, 33.0, \dots]$ باشد. این نمایش باعث می‌شود که مدل بتواند ارتباطات معنایی بین کلمات را بهتر درک کند و یاد بگیرد که cat و dog از نظر معنایی به هم نزدیک‌تر از cat و banana هستند. این ویژگی در مدل‌های توضیح تصویر اهمیت زیادی دارد زیرا مدل باید بتواند از ویژگی‌های تصویری، کلماتی را انتخاب کند که به طور معنایی به تصویر مرتبط باشند.

یکی دیگر از مزایای مهم Word Embedding این است که باعث کاهش ابعاد و بهینه‌سازی محاسبات می‌شود. همان‌طور که گفتیم در روش One-Hot Encoding بردارهای نمایش داده شده برای کلمات بسیار طولانی هستند و اگر واژگان مدل شامل دهها هزار کلمه باشد، هر کلمه باید در برداری با همین تعداد بعد نمایش داده شود که از نظر محاسباتی بسیار سنگین و ناکارآمد است. اما در Word Embedding هر کلمه در یک فضای کم‌بعد نمایش داده می‌شود، مثلاً در فضایی با ۲۵۶ بعد که باعث کاهش مصرف حافظه و افزایش کارایی مدل می‌شود. این کار به مدل کمک می‌کند که سریع‌تر آموزش ببیند و بتواند مجموعه داده‌های بزرگ را بدون نیاز به منابع زیاد پردازش کند.

علاوه بر این، استفاده از Word Embedding باعث بهبود تعمیم‌پذیری مدل می‌شود. در روش One-Hot Encoding، هر کلمه به صورت مجزا نمایش داده می‌شود و اگر یک کلمه جدید در داده‌های آزمایشی

ظاهر شود که در زمان آموزش وجود نداشته است، مدل قادر به پردازش آن نخواهد بود زیرا هیچ برداری برای آن در نظر گرفته نشده است. اما در Word Embedding، حتی اگر یک کلمه جدید باشد، مدل می‌تواند بر اساس بردارهای کلمات مشابه، معنای آن را تا حدی تخمین بزند و همچنان عملکرد مطلوبی داشته باشد. این ویژگی بهویژه در مدل‌های توضیح تصویر مفید است زیرا تصاویر جدید ممکن است شامل کلماتی باشند که قبلًا در مجموعه داده آموزش وجود نداشتند.

مزیت دیگری که Word Embedding دارد این است که می‌توان از مدل‌های از پیشآموزش داده شده مانند FastText یا Word2Vec استفاده کرد. این مدل‌ها روی مجموعه داده‌های بسیار بزرگ مانند Google News و Wikipedia آموزش دیده‌اند و بردارهای بسیار معناداری برای کلمات ایجاد کرده‌اند. در نتیجه، به جای یادگیری بردارهای جدید از ابتدا، می‌توان از بردارهای از پیشآموزش داده شده استفاده کرد که هم دقت مدل را بالا می‌برد و هم زمان آموزش را کاهش می‌دهد. این کار باعث می‌شود که مدل ما نیازی به یادگیری مجدد معنای کلمات نداشته باشد و بتواند از دانش قبلی که در مدل‌های Word Embedding ذخیره شده است استفاده کند.

در پیاده‌سازی مدل توضیح تصویر، ما یک لایه تعبیه کلمات را اضافه کردیم تا کلمات ورودی ابتدا به بردارهای یادگیری شده تبدیل شوند. این لایه مسئول تبدیل هر کلمه به یک بردار عددی معنادار است که سپس به LSTM داده می‌شود تا پردازش ترتیب کلمات را انجام دهد. در واقع، لایه Embedding نقش یک نگاشت یادگیری شده را ایفا می‌کند که هر کلمه را به یک نقطه در فضای برداری انتقال می‌دهد. این باعث می‌شود که مدل بتواند اطلاعات معنایی کلمات را حفظ کند و آن‌ها را در کنار یکدیگر در جمله به شکلی معنادار استفاده کند. در نهایت، این نمایش برداری کمک می‌کند که مدل بتواند کلمات را به درستی تشخیص داده و ارتباط بین آن‌ها را درک کند که در نهایت منجر به تولید کپشن‌های روان‌تر و معنادار‌تر می‌شود. به همین دلیل، استفاده از Word Embedding در این پروژه یکی از اجزای کلیدی مدل ما بوده است و کمک زیادی به بهبود عملکرد آن کرده است.

چگونه رمزگذار و رمزگشا را به یک مدل End-to-End تبدیل کنیم که قابلیت آموزش داشته باشد؟

برای تبدیل رمزگذار و رمزگشا به یک مدل End-to-End که قابلیت آموزش داشته باشد، باید این دو بخش را در قالب یک معماری یکپارچه طراحی کنیم که در آن رمزگذار (Encoder) ویژگی‌های تصویر را استخراج کند و رمزگشا (Decoder) از این ویژگی‌ها برای تولید متن توصیفی استفاده کند. در واقع، این فرآیند شامل ترکیب یک شبکه عصبی پیچشی (CNN) برای استخراج ویژگی‌های تصویری و یک شبکه عصبی بازگشتی (RNN) یا مدل‌های پیشرفته‌تر مانند LSTM یا GRU برای تولید جملات کپشن است.

به این ترتیب، ما یک مدل یادگیری عمیق را ایجاد می کنیم که می تواند مستقیماً از تصاویر ورودی، خروجی متنی تولید کند. برای این کار، چندین گام باید انجام شود که در ادامه توضیح خواهیم داد.

در مرحله‌ی اول، باید رمزگذار را طراحی کنیم که در اینجا از مدل EfficientNet-B0 استفاده کردہ‌ایم. رمزگذار وظیفه دارد تصویر ورودی را پردازش کرده و آن را به یک بردار ویژگی با ابعاد فشرده تبدیل کند که نشان‌دهنده‌ی محتوای معنایی تصویر است. مدل‌های CNN مانند EfficientNet برای این کار بهینه شده‌اند زیرا قادرند ویژگی‌های سطح پایین و سطح بالای تصویر را استخراج کنند و اطلاعات مفیدی از محتویات تصویر به دست آورند. در این پیاده‌سازی، ما آخرین لایه Fully Connected مدل EfficientNet-B0 را حذف کردیم تا تنها ویژگی‌های خام تصویر استخراج شوند. سپس این ویژگی‌ها را از طریق یک لایه خطی به بعدی کاهش دادیم که برای رمزگشا قابل استفاده باشد. این مرحله بسیار مهم است، زیرا خروجی CNN دارای تعداد زیادی ویژگی‌های فضایی است که باید به یک بردار فشرده تبدیل شوند تا قابل استفاده در LSTM باشند. بعد از این مرحله، رمزگذار آماده است که ویژگی‌های تصویری را استخراج کرده و به رمزگشا ارسال کند.

در مرحله‌ی دوم، باید رمزگشا را طراحی کنیم که بتواند از بردار ویژگی استخراج شده توسط رمزگذار برای تولید توصیف تصویر استفاده کند. در اینجا از یک مدل LSTM به عنوان رمزگشا استفاده کردہ‌ایم که وظیفه‌ی آن تولید توالی کلمات است. برای این کار، ابتدا یک لایه تعبیه کلمات (Word Embedding) در رمزگشا قرار داده‌ایم که مسئول تبدیل کلمات ورودی به بردارهای معنایی فشرده است. این کار به مدل کمک می‌کند تا روابط معنایی بین کلمات را درک کند و کمیشن‌های معنادارتری تولید کند. سپس این بردارهای کلمات همراه با ویژگی‌های تصویری استخراج شده توسط رمزگذار وارد LSTM می‌شوند. نکته‌ی مهم این است که ویژگی‌های تصویری باید به عنوان حالت اولیه (Initial Hidden State) به LSTM داده شوند، زیرا LSTM یک مدل ترتیبی است که نیاز دارد از جایی پردازش خود را شروع کند. برای این کار، ما بردار ویژگی‌های تصویر را به عنوان ورودی اولیه به LSTM می‌دهیم و سپس مدل شروع به پردازش دنباله‌ی کلمات می‌کند.

پس از گذراندن بردارهای تعبیه شده از طریق LSTM، مدل خروجی‌های مربوط به کلمات را تولید می‌کند، اما این خروجی‌ها به صورت بردارهای عددی هستند که باید به کلمات قابل فهم تبدیل شوند. برای این کار، یک لایه خطی (Fully Connected Layer) همراه با Softmax را در انتهای LSTM قرار داده‌ایم که احتمال هر کلمه را در مرحله‌ی بعدی محاسبه می‌کند و در نهایت کلمه‌ای که بالاترین احتمال را دارد انتخاب می‌شود. این کار به مدل اجازه می‌دهد که برای هر مرحله، کلمه‌ای مناسب را از واژگان انتخاب کند و جمله‌ی نهایی را تشکیل دهد.

برای این که رمزگذار و رمزگشا را در قالب یک مدل واحد که قابلیت آموزش داشته باشد پیاده‌سازی کنیم، آن‌ها را درون یک کلاس سفارشی به نام ImageCaptioningModel قرار داده‌ایم که شامل دو بخش Encoder و Decoder است. این کلاس ابتدا تصویر را به رمزگذار داده و ویژگی‌های استخراج شده را دریافت می‌کند. سپس این ویژگی‌ها به رمزگشا داده می‌شوند تا متن کپشن تولید شود. یکی از نکات کلیدی در طراحی این مدل این است که باید امکان انتشار گرادیان‌ها از رمزگشا به رمزگذار فراهم باشد تا مدل بتواند به صورت End-to-End آموزش ببینند. برای این کار، اتصال بین این دو بخش باید به گونه‌ای باشد که هنگام به روزرسانی وزن‌ها در طول آموزش، گرادیان‌ها از رمزگشا به رمزگذار منتقل شوند.

برای آموزش این مدل، ابتدا یک تابع هزینه مناسب مانند CrossEntropyLoss را تعریف کردیم که برای مدل‌های تولید متن به خوبی عمل می‌کند. همچنین برای این که مدل بیش از حد روی نمونه‌های خاص وابسته نشود، Label Smoothing را به مقدار $1 - \epsilon$ اضافه کردیم. یکی دیگر از نکات مهم در آموزش این مدل، نادیده گرفتن Padding در محاسبه‌ی هزینه است، زیرا در جملات با طول‌های متغیر، برخی از کلمات اضافه‌شده تنها برای پر کردن توالی هستند و نباید در محاسبات تأثیر بگذارند. سپس برای بهبود همگرایی مدل کمک می‌کند. همچنین از AdamW استفاده کردیم که عملکرد مناسبی برای مدل‌های یادگیری عمیق دارد و به بهبود همگرایی مدل ReduceLROnPlateau برای کاهش نرخ یادگیری در صورت عدم پیشرفت استفاده کردیم.

در طول آموزش، ابتدا تصاویر و کپشن‌های مربوطه را از دیتابست دریافت کرده و تصاویر را از طریق رمزگذار پردازش کردیم تا ویژگی‌های آن‌ها استخراج شوند. سپس این ویژگی‌ها همراه با توالی کلمات وارد رمزگشا شدند تا مدل بتواند توصیف تصویر را مرحله‌به‌مرحله تولید کند. در طول هر مرحله‌ی آموزش، مدل با مقایسه‌ی خروجی‌های پیش‌بینی‌شده با کپشن‌های واقعی، خطای خود را محاسبه کرد و وزن‌های خود را بهروز کرد تا به تدریج دقت خود را افزایش دهد. همچنین، در پایان هر Epoch، نمونه‌ای از تصاویر و کپشن‌های تولیدشده را نمایش دادیم تا عملکرد مدل را بررسی کنیم. برای ارزیابی نهایی، از روش Beam Search برای تولید توصیف‌های بهتر استفاده کردیم. این روش به جای این که در هر مرحله فقط محتمل‌ترین کلمه را انتخاب کند، چندین کلمه‌ی برتر را بررسی می‌کند و بهترین ترکیب را انتخاب می‌کند تا خروجی نهایی روان‌تر و دقیق‌تر باشد.

در نهایت، پس از آموزش مدل، پنج تصویر به همراه کپشن‌های تولیدشده را نمایش دادیم تا بررسی کنیم که مدل چگونه محتوای تصاویر را توصیف کرده است. مشاهده شد که مدل توانسته است بسیاری از اشیاء و فعالیت‌های موجود در تصاویر را به درستی تشخیص دهد و توصیف کند، اما در برخی موارد همچنان خطاطهایی مانند عدم تشخیص دقیق برخی اشیاء یا اشتباه در ارتباطات بین آن‌ها مشاهده شد. برای بهبود

بیشتر مدل، می‌توان در آینده از مکانیزم Attention استفاده کرد که به رمزگشا اجازه می‌دهد روی بخش‌های خاصی از تصویر تمرکز کند و اطلاعات مهم‌تر را در هنگام تولید کلمات مدنظر قرار دهد. در نهایت، مدل آموزش‌یافته ذخیره شد و اکنون می‌توان از آن برای تولید کپشن‌های جدید بر روی تصاویر استفاده کرد.

بررسی کنید مقاله چه تکنیک‌هایی برای جلوگیری از بیشبرازش Overfitting () را بکار می‌گیرد؟

مقاله "Show and Tell: A Neural Image Caption Generator" برای جلوگیری از بیشبرازش در مدل توضیح تصویر، از چندین تکنیک مختلف استفاده کرده است. در این بخش، Overfitting) تکنیک‌های به کار رفته در این مقاله برای کاهش بیشبرازش و بهبود تعمیم‌پذیری مدل بررسی می‌شود. مدل‌های یادگیری عمیق، به خصوص آن‌هایی که روی مجموعه داده‌های نسبتاً کوچک آموزش داده می‌شوند، به راحتی دچار بیشبرازش می‌شوند. در این مقاله، به دلیل آنکه داده‌های تصویری حاوی اطلاعات پیچیده‌ای هستند و مدل نیاز دارد تا جملات معنایی تولید کند، مسئله بیشبرازش چالش بسیار مهمی است. نویسنده‌گان مقاله به این موضوع توجه ویژه‌ای داشته‌اند و چندین روش برای کاهش بیشبرازش و افزایش کارایی مدل در داده‌های جدید ارائه کرده‌اند.

یکی از مهم‌ترین تکنیک‌های استفاده شده، انتقال یادگیری (Transfer Learning) است. در این روش، نویسنده‌گان مقاله برای جلوگیری از بیشبرازش، وزن‌های شبکه عصبی پیچشی CNN را با استفاده از یک مدل از پیشآموزش داده شده (Pretrained Model) مقداردهی اولیه کرده‌اند. این مدل از پیشآموزش داده شده روی دیتابست بزرگ ImageNet آموزش دیده است و ویژگی‌های تصویری را به خوبی استخراج می‌کند. استفاده از مدل از پیشآموزش داده شده باعث می‌شود که مدل نیازی به یادگیری از صفر نداشته باشد و از دانشی که قبلاً در شبکه یاد گرفته شده، استفاده کند. این کار باعث می‌شود که مدل سریع‌تر همگرا شود و از یادگیری بیش از حد بر روی داده‌های آموزشی جلوگیری شود. به همین دلیل، در تمام آزمایش‌های انجام شده در مقاله، وزن‌های شبکه CNN با مقداردهی اولیه از مدل آموزش داده شده روی ImageNet استفاده شده است.

تکنیک دیگر، استفاده از Dropout در لایه‌های مختلف مدل است. Dropout یکی از رایج‌ترین روش‌های جلوگیری از بیشبرازش در شبکه‌های عصبی است که با حذف تصادفی برخی از نورون‌ها در طول آموزش، از وابستگی بیش از حد مدل به مجموعه داده‌ی آموزشی جلوگیری می‌کند. در این مقاله، Dropout به عنوان یک روش مؤثر برای بهبود تعمیم‌پذیری مدل مورد استفاده قرار گرفته است. این تکنیک باعث می‌شود که مدل به جای تکیه بر ویژگی‌های خاص در تصاویر آموزشی، الگوهای کلی‌تر و عمومی‌تری یاد بگیرد که در

مواجهه با داده‌های جدید، عملکرد بهتری داشته باشد. نویسنده‌گان مقاله گزارش کرده‌اند که استفاده از Dropout چندین امتیاز BLEU در بهبود عملکرد مدل اضافه کرده است.

علاوه بر Dropout، نویسنده‌گان از Ensemble Learning نیز استفاده کرده‌اند. در این روش، چندین مدل مختلف آموزش داده شده‌اند و سپس خروجی آن‌ها با هم ترکیب شده تا پیش‌بینی نهایی بهبود یابد. استفاده از Ensemble باعث می‌شود که خطاهای تصادفی که ممکن است در یک مدل منفرد وجود داشته باشد، کاهش یابد و مدل کلی به تعمیم‌پذیری بهتری برسد. در مقاله ذکر شده که این تکنیک نیز به بهبود چندین امتیاز BLEU در مدل نهایی کمک کرده است.

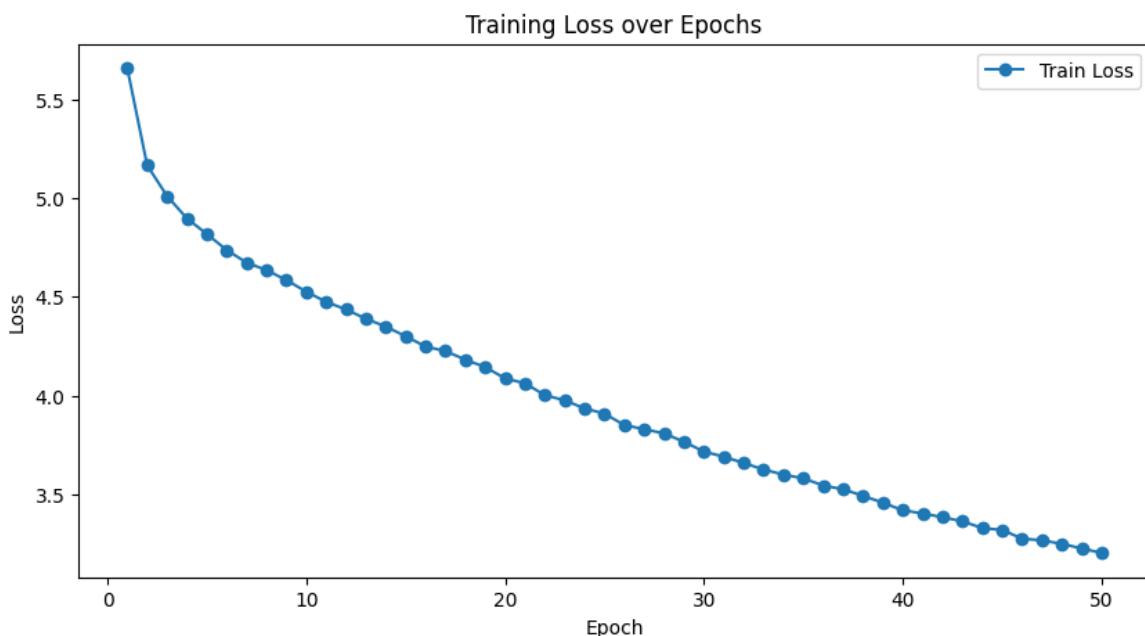
یکی دیگر از روش‌هایی که در مقاله برای کاهش بیش‌برازش به کار گرفته شده، تنظیم ابعاد مدل (Capacity Control) است. در این روش، نویسنده‌گان اندازه مدل را با تغییر تعداد نورون‌های مخفی در لایه‌های LSTM تنظیم کرده‌اند و بین عمق شبکه و تعداد واحدهای مخفی تعادل ایجاد کرده‌اند. اگر مدل بیش‌از حد پیچیده باشد، احتمال دارد که بر روی داده‌های آموزشی بیش‌برازش کند و در مواجهه با داده‌های جدید عملکرد خوبی نداشته باشد. از سوی دیگر، اگر مدل بیش‌از حد کوچک باشد، ممکن است قدرت کافی برای یادگیری ویژگی‌های لازم را نداشته باشد. در این مقاله، آزمایش‌های متعددی برای یافتن بهترین اندازه مدل انجام شده است تا از بیش‌برازش جلوگیری شود و در عین حال مدل قدرت کافی برای یادگیری داشته باشد.

نویسنده‌گان همچنین نرمال‌سازی دسته‌ای (Batch Normalization) به کار گرفته‌اند. این روش باعث می‌شود که مقادیر ورودی هر لایه در طول آموزش ثابت نگه داشته شوند و از تغییرات ناگهانی در گرادیان جلوگیری شود. Batch Normalization نه تنها باعث سرعت‌بخشیدن به فرآیند آموزش می‌شود، بلکه به جلوگیری از بیش‌برازش نیز کمک می‌کند، زیرا باعث می‌شود که مدل با تغییرات ناگهانی در داده‌های آموزشی بیش‌از حد تنظیم نشود و تعمیم‌پذیری بهتری داشته باشد.

نویسنده‌گان مقاله در بخشی دیگر تلاش کرده‌اند که واژگان مدل (Vocabulary) را کنترل کنند تا از بیش‌برازش جلوگیری کنند. در این راستا، آن‌ها فقط کلماتی را در واژگان نگه داشته‌اند که حداقل پنج بار در مجموعه داده آموزشی تکرار شده باشند. این کار باعث می‌شود که مدل بر روی کلمات پر تکرارتر و مهم‌تر تمرکز کند و از یادگیری نویزهای موجود در داده‌ها جلوگیری کند. همچنین، توکن‌های ویژه‌ای مانند <eos> و <sos> برای مشخص کردن آغاز و پایان جملات استفاده شده‌اند که به مدل کمک می‌کند توصیف‌های منسجم‌تری تولید کند.

برای بهینه‌سازی فرآیند آموزش، نویسنده‌گان از روش کاهش نرخ یادگیری (Learning Rate Scheduling) نیز استفاده کردند. در این روش، اگر مدل برای چندین epoch متواتی نتواند بهبود قابل توجهی در معیار ارزیابی نشان دهد، نرخ یادگیری کاهش پیدا می‌کند تا مدل بتواند به شکلی پایدارتر و دقیق‌تر بهینه شود. این روش کمک می‌کند که مدل در اوایل آموزش سریع یاد بگیرد و در مراحل بعدی از نوسانات زیاد جلوگیری شود.

نمودار خطای داده آموزش را در طول هر دوره Epoch گزارش کنید.



شکل ۸ نمودار خطای آموزش در طول هر دوره

این نمودار نشان می‌دهد که مقدار خطای آموزش (Loss) در ابتدای آموزش نسبتاً بالا بوده و در طی ۵۰ دوره به تدریج کاهش یافته است. این روند نزولی نشانه‌ی مثبتی است و بیانگر این است که مدل به تدریج در حال یادگیری الگوهای موجود در داده‌های آموزشی است. در ابتدای آموزش، خطا در حدود ۵.۶ بوده و در طول دوره‌های بعدی به کمتر از ۳.۰ کاهش یافته است. این نشان می‌دهد که مدل به تدریج در حال کاهش خطای خود و افزایش دقت در پیش‌بینی کلمات مناسب برای کپشن‌ها است.

یکی از نکات قابل مشاهده در این نمودار این است که کاهش خطا در چند epoch اول بسیار سریع است، اما هرچه تعداد epoch‌ها افزایش پیدا می‌کند، نرخ کاهش خطا کندتر می‌شود. این رفتار طبیعی است، زیرا در ابتدای آموزش، مدل سریع‌ترین میزان یادگیری را دارد و هرچه بیشتر یاد می‌گیرد، بهبودهای بعدی دشوارتر می‌شوند.

اگرچه این نمودار نشان‌دهنده‌ی روند کاهش خطای آموزش است، اما برای اطمینان از این که مدل دچار بیش‌برازش (Overfitting) نشده، نیاز به مقایسه‌ی آن با نمودار خطای اعتبارسنجی (Validation Loss) داریم. اگر در این نمودار مشاهده می‌شد که پس از یک نقطه خاص، خطای آموزش همچنان کاهش می‌یابد اما خطای اعتبارسنجی افزایش پیدا می‌کند، نشان‌دهنده‌ی بیش‌برازش بود. از آنجایی که این نمودار فقط خطای آموزش را نمایش می‌دهد، نمی‌توان از این موضوع مطمئن شد، اما اگر خطای آموزش همچنان در حال کاهش باشد و خطای اعتبارسنجی نیز کاهش پیدا کند، می‌توان نتیجه گرفت که مدل به درستی در حال یادگیری است.

با این حال، از شکل کلی این نمودار می‌توان حدس زد که هنوز مدل در مرحله‌ای است که به درستی در حال یادگیری تعمیم یافته (Generalized Learning) است. اگر روند کاهش خطا در انتهای نمودار خیلی زود متوقف می‌شود یا حتی شروع به افزایش می‌کرد، این احتمال وجود داشت که مدل در حال بیش‌برازش باشد. اما از آنجا که خطا همچنان در حال کاهش است، نشانه‌ی مثبتی از عدم بیش‌برازش مشاهده می‌شود.

یکی از جنبه‌های مهمی که می‌توان از این نمودار استخراج کرد، نرخ یادگیری مدل (Learning Rate) است. از شکل نمودار مشخص است که نرخ یادگیری در طول آموزش مناسب بوده است. اگر نرخ یادگیری بیش از حد بالا تنظیم می‌شود، مدل ممکن بود در نوسانات شدید گیر بیفتند و خطای آموزش به درستی کاهش پیدا نمی‌کرد. همچنین، اگر نرخ یادگیری بیش از حد کم بود، کاهش خطا بسیار کندتر اتفاق می‌افتد و مدل زمان زیادی برای یادگیری نیاز داشت. روند هموار و تدریجی کاهش خطا نشان می‌دهد که مقدار نرخ یادگیری تنظیم شده است.

همچنین، در مراحل پایانی آموزش، روند کاهش خطا همچنان ادامه دارد، اما شیب آن نسبت به مراحل ابتدایی کمتر شده است. این نشان می‌دهد که مدل در حال همگرایی (Convergence) است و به نقطه‌ای می‌رسد که یادگیری به حداقل خود رسیده است. معمولاً در چنین مواردی، می‌توان از تکنیک‌های کاهش نرخ یادگیری دینامیکی (Learning Rate Scheduling) استفاده کرد که باعث می‌شود نرخ یادگیری به مرور زمان کاهش یابد و مدل با دقت بیشتری یادگیری کند.

یکی از روش‌های جلوگیری از آموزش بیش از حد (Overtraining) استفاده از توقف زودهنگام (Early Stopping) است. در این روش، اگر مشاهده شود که خطای آموزش همچنان کاهش می‌یابد اما خطای اعتبارسنجی بهبودی ندارد یا افزایش می‌یابد، می‌توان آموزش را متوقف کرد تا از بیش‌برازش جلوگیری شود. در این نمودار، کاهش خطا همچنان ادامه دارد، اما در اواخر آموزش سرعت آن کمتر شده است. بررسی نمودار خطای اعتبارسنجی می‌تواند نشان دهد که آیا در اواخر آموزش نیاز به استفاده از توقف زودهنگام وجود دارد یا خیر.

کمآموزی (Underfitting) زمانی رخ می‌دهد که مدل هنوز نتوانسته ویژگی‌های لازم را از داده‌ها استخراج کند و خطای آموزش در مقدار بالایی متوقف شود. در این نمودار، خطای آموزش از مقدار ۵,۶ به کمتر از ۳,۰ کاهش یافته است که نشان می‌دهد مدل توانسته است مقدار قابل توجهی از ویژگی‌های داده را یاد بگیرد. اگر مقدار خطای آموزش در حدود ۴ یا ۵ باقی می‌ماند و کاهش نمی‌یابد، این احتمال وجود داشت که مدل دچار کمآموزی باشد. اما در اینجا، روند نزولی پیوسته نشان‌دهنده‌ی یادگیری موثر مدل است.

میتوانید در پایان هر دوره یک نمونه تصویر و خروجی آن را نمایش دهید.

این کار در کد پروژه انجام شده است به دلیل تعدد تصاویر تمامی عکس‌ها در اینجا قرار داده نمی‌شود صرفاً اولین تصویر خروجی در اولین دوره اموزشی در زیر قرار داده شده است.

Pred: three man on sidewalk on sidewalk with red is on sidewalk with ball with ball
GT: two guys dressed in tuxes with red ties



شکل ۹ پیش‌بینی مدل در دوره اول

```
Epoch 1/50:  0% | 1/203 [00:00<00:30,  6.72it/s, loss=9.09]
[INFO] Encoder output shape: torch.Size([32, 256])
[INFO] Decoder output shape: torch.Size([32, 15, 8781])
Epoch 1/50: 100%|██████████| 203/203 [00:30<00:00,  6.61it/s, loss=5.75]
[INFO] Checkpoint saved at epoch 0
Epoch 1/50, Loss: 5.6610
```

در پایان آموزش، ۵تا از تصاویر و توضیحات تولیدشده آنرا در کنار یکدیگر نمایش دهید.
برخی از خطاهای مدل را شناسایی و مشخص کنید (مانند عدم تشخیص اشیاء یا روابط)

در این بخش از پروژه، به ارزیابی مدل توضیح تصویر (Image Captioning) پرداختیم تا عملکرد آن را در تولید کپشن‌های توصیفی برای تصاویر بسنجیم. هدف اصلی این ارزیابی، بررسی کیفیت کپشن‌های تولیدشده توسط مدل است و از معیارهای BLEU و ROUGE برای این منظور استفاده کردیم. معیارهای BLEU و ROUGE به ما کمک می‌کنند تا بفهمیم مدل تا چه اندازه توانسته کپشن‌هایی نزدیک به کپشن‌های واقعی و انسانی تولید کند. این معیارها با مقایسه توالی‌های کلمات در کپشن‌های تولیدی و کپشن‌های مرجع، کیفیت و دقت مدل را مورد سنجش قرار می‌دهند.

ابتدا مدل آموزش‌دهنده را که قبلًا ذخیره کرده بودیم، بارگذاری کردیم. این کار با استفاده از تابع torch.load انجام شد و مدل به دستگاه محاسباتی GPU یا (CPU) یا انتقال داده شد. همچنین، مدل را در حالت ارزیابی قرار دادیم تا اطمینان حاصل کنیم که در طول ارزیابی هیچ تغییری در وزن‌های شبکه ایجاد نمی‌شود. این کار با استفاده از model.eval() انجام شد که حالت آموزش را غیرفعال می‌کند و مدل فقط برای پیش‌بینی مورد استفاده قرار می‌گیرد.

برای ارزیابی کیفیت کپشن‌های تولیدی از معیار BLEU استفاده کردیم که یکی از شناخته‌شده‌ترین روش‌ها برای ارزیابی کیفیت ترجمه‌های ماشینی و تولید متون توسط مدل‌های یادگیری عمیق است. در این پروژه، چهار نسخه از نمره BLEU را محاسبه کردیم که هر یک به بررسی میزان شباهت کپشن‌های تولیدی با کپشن‌های مرجع در سطوح‌های مختلف n-gram می‌پردازند. BLEU-1 تطابق تک کلمه‌ای (unigram) را می‌سنجد، BLEU-2 به توالی‌های دو کلمه‌ای (bigram) می‌پردازد، BLEU-3 به توالی‌های سه کلمه‌ای (trigram) می‌پردازد، BLEU-4 به توالی‌های چهار کلمه‌ای (4-gram) استفاده می‌شود. برای محاسبه این نمرات، از تابع corpus_bleu از کتابخانه NLTK استفاده کردیم که به کمک روش هموارسازی Smoothing عملکرد بهتری روی جملات کوتاه دارد. این روش به ویژه زمانی اهمیت دارد که مدل جملات کوتاهی تولید می‌کند که ممکن است به دلیل کمبود داده یا محدودیت‌های ساختاری، برخی کلمات ضروری را شامل نشوند.

علاوه بر BLEU، از معیار ROUGE نیز استفاده کردیم که بیشتر برای ارزیابی خلاصه‌سازی متن به کار می‌رود، اما در اینجا برای بررسی کیفیت تولید کپشن‌های تصویر استفاده شد. ما از نسخه ROUGE-L استفاده کردیم که بر پایه‌ی بزرگترین دنباله مشترک (Longest Common Subsequence) بین کپشن مرجع و کپشن تولیدی عمل می‌کند. این معیار به ما کمک می‌کند تا بفهمیم چه مقدار از کلمات یا عبارات مهم در کپشن‌های تولیدی با کپشن‌های مرجع همپوشانی دارند. استفاده از ROUGE-L به مدل اجازه می‌دهد تا نه تنها بر اساس توالی کلمات، بلکه بر اساس اهمیت معنایی آن‌ها نیز ارزیابی شود.

در بخش بعدی، نیاز داشتیم تا خروجی مدل را به صورت متن قابل خواندن تبدیل کنیم. این کار با استفاده از تابع `clean_caption` انجام شد. این تابع ابتدا شاخص‌های عددی که مدل تولید کرده بود را به کلمات واقعی تبدیل می‌کند. همچنین، توکن‌های اضافی مانند `<sos>` و `<eos>` و `<pad>` که تنها برای اهداف فنی در زمان آموزش استفاده می‌شوند را حذف کردیم تا کپشن‌های نهایی به صورت جملات روان و قابل فهم باشند. این تابع علاوه بر حذف توکن‌های اضافی، کلمات تکراری پشت سر هم را نیز حذف می‌کند تا از تولید جملات نامفهوم جلوگیری شود.

برای اجرای فرآیند ارزیابی، تابع `evaluate_model` را تعریف کردیم که مدل را روی داده‌های مختلف اعتبارسنجی و تست ارزیابی می‌کند. در این تابع، ابتدا مدل به حالت ارزیابی درآمده و سپس تصاویر و کپشن‌های مرجع از داده‌بار (DataLoader) دریافت می‌شوند. برای هر تصویر، مدل خروجی خود را به صورت توکن‌های پیش‌بینی‌شده تولید می‌کند. این توکن‌ها سپس به کپشن‌های متنی تبدیل می‌شوند تا با کپشن‌های مرجع مقایسه شوند. در نهایت، نمرات BLEU و ROUGE برای کل مجموعه داده محاسبه و نمایش داده می‌شوند. این نمرات به ما نشان می‌دهند که مدل چقدر توانسته به درستی محتواهای تصاویر را توصیف کند و کپشن‌هایی مشابه با کپشن‌های انسانی تولید نماید.

بهمنظور درک بهتر عملکرد مدل، پنج نمونه تصویر به صورت تصادفی انتخاب شده و کپشن‌های تولیدی و مرجع آن‌ها به همراه تصویر نمایش داده شدند. این نمایش بصری به ما کمک می‌کند که به‌طور کیفی نیز بررسی کنیم که آیا مدل توانسته ارتباطات معنایی و بصری موجود در تصاویر را درک کند یا خیر. به عنوان مثال، اگر مدل بتواند به درستی تشخیص دهد که یک سگ در حال بازی است یا یک کودک در حال خنده است، این نشان می‌دهد که مدل به‌خوبی یاد گرفته است که ویژگی‌های تصویری را به متن مرتبط تبدیل کند. اما اگر مدل به اشتباه گربه را به جای سگ تشخیص دهد یا فعالیت‌ها را به درستی توصیف نکند، نشان می‌دهد که هنوز نیاز به بهبود در برخی از جنبه‌ها وجود دارد.

Real: a brown dog on a leash runs through the white water
Pred: a brown dog is a beach in the water



Real: a group of people sit around a table outside on a porch at night
Pred: a group of people are in a campfire in a brick day a



Real: a young brown haired woman plays her guitar and sings into a microphone
Pred: a man in wearing in a violin in a microphone



Real: the dog is running along a path that has been made through the uncut
Pred: a brown dog running through the dirt



Real: a man on a bmx bike jumps over a train
Pred: a skateboarder in a bicycle is performing midair a chain

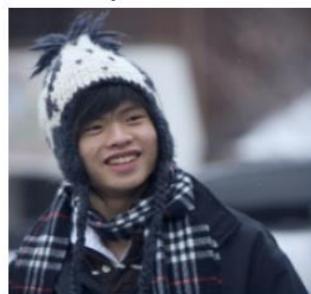


شکل ۱۰ نمونه هایی از کپشن های ایجاد شده در داده ولیدیشن

Real: a kid with sunglasses in jumping on the beach
Pred: a young boy a bucket on white a beach



Real: a man wearing a scarf and knitted hat is smiling
Pred: a woman boy a and shirt down hat looks a



Real: two women walk down a city street at night
Pred: a man in a street



Real: two children one of which is holding a stick are standing on the beach
Pred: two children are in a s on a and they towards a shoulders



Real: a dog swims in the aqua water
Pred: a dog in a pool



شکل ۱۱ نمونه هایی از کپشن های ایجاد شده در داده تست

در بررسی خروجی‌های مدل روی داده‌های تست و اعتبارسنجی، مشاهده می‌شود که مدل در بسیاری از موارد توانسته است مفاهیم کلی تصاویر را درک کند، اما همچنان با مشکلاتی در دقت و جزئیات توصیف‌ها مواجه است. برخی از کپشن‌های تولیدشده کاملاً نزدیک به کپشن‌های واقعی هستند، اما در برخی نمونه‌ها، مدل دچار اشتباهاتی در تشخیص سوژه‌ها و ارتباطات بین آن‌ها شده است. یکی از دلایل اصلی این خطاهای محدودیت‌های ذاتی مدل‌های RNN و LSTM در درک وابستگی‌های طولانی مدت در داده‌های متوالی است.

یکی از مهم‌ترین مشکلات مشاهده شده، عدم دقت در تشخیص دقیق سوژه‌ها است. برای مثال، در تصویر اول، کپشن واقعی به این نکته اشاره می‌کند که یک سگ قهوه‌ای با قلاده در حال دویدن در میان آب سفید است، اما مدل این صحنه را به صورت "یک سگ قهوه‌ای در ساحل" توصیف کرده است. این نشان می‌دهد که مدل مفهوم کلی را درک کرده، اما قادر به تشخیص جزئیاتی مانند "دویدن" و "قلاده" نبوده است. این مسئله در بسیاری از نمونه‌ها دیده می‌شود که مدل می‌تواند سوژه‌ی اصلی تصویر را شناسایی کند، اما در تشخیص زمینه‌ی تصویر و فعالیت‌های مرتبط دچار مشکل می‌شود.

مشکل دیگری که مشاهده می‌شود، عدم درک دقیق مفاهیم انسانی و فعالیت‌ها است. برای مثال، در تصویری که کپشن واقعی آن "زنی با موهای قهوه‌ای در حال نواختن گیتار و خواندن در میکروفون" است، مدل این تصویر را به عنوان "مردی که ویولون می‌نوازد" توصیف کرده است. این نشان می‌دهد که مدل در تمایز بین جنسیت افراد و همچنین نوع فعالیت‌های مرتبط با موسیقی دچار ضعف است. این نوع خطاهای معمولاً به دلیل محدودیت‌های RNN در پردازش همزمان اطلاعات بصری و زبانی ایجاد می‌شود، زیرا این مدل‌ها تنها اطلاعات گذشته را به حافظه می‌سپارند و ارتباطات طولانی مدت در توالی داده‌ها را به خوبی حفظ نمی‌کنند.

مشکل دیگری که در برخی نمونه‌ها به چشم می‌خورد، عدم حفظ ساختار جملات و روان نبودن برخی خروجی‌ها است. در بعضی از کپشن‌های تولیدشده، کلمات اضافی، ساختارهای ناهماهنگ و ترکیب‌های نامفهوم مشاهده می‌شود. برای مثال، در یکی از نمونه‌ها که مدل سعی کرده است یک صحنه‌ی گروهی را توصیف کند، به جای جمله‌ای روان و معنادار، عبارتی غیرمنسجم تولید کرده که مفهوم روشنی ندارد. این مشکل نیز به ساختار ترتیبی RNN برمی‌گردد، زیرا این مدل‌ها جمله را به صورت گام‌به‌گام تولید می‌کنند و نمی‌توانند ارتباطات طولانی مدت و پیچیدگی‌های نحوی زبان را به خوبی مدیریت کنند.

یکی دیگر از نقاط ضعف مدل استفاده شده، عدم تمرکز روی جزئیات مهم تصویر است. مدل‌های RNN به دلیل پردازش ترتیبی داده‌ها، نمی‌توانند همزمان روی بخش‌های مختلف تصویر تمرکز کنند و این موضوع باعث می‌شود که اطلاعات مهم تصویری در حین تولید کپشن از دست برود. در برخی نمونه‌ها، مدل قادر

به تشخیص سوژه‌ی اصلی تصویر بوده، اما در توصیف موقعیت مکانی، پس‌زمینه و جزئیات دقیق دچار مشکل شده است. این مسئله دقیقاً همان جایی است که مدل‌های Attention و ترانسفورمرها می‌توانند به شکل قابل توجهی عملکرد را بهبود دهند.

در نسخه‌های جدیدتر مدل‌های توضیح تصویر، از مکانیزم توجه (Attention Mechanism) استفاده می‌شود که به مدل اجازه می‌دهد تا در هر لحظه از تولید کپشن، بر روی بخش‌های خاصی از تصویر تمرکز کند. به این ترتیب، مدل قادر خواهد بود نه تنها اشیاء را بهتر شناسایی کند، بلکه موقعیت مکانی آن‌ها و ارتباطات میان آن‌ها را نیز در نظر بگیرد. علاوه بر این، در مدل‌های مدرن‌تر مانند ترانسفورمرها (Transformers) که در معماری‌هایی مانند Vision Transformer (ViT) و BLIP به کار گرفته شده‌اند، مدل می‌تواند هم‌زمان تمام قسمت‌های تصویر و کپشن را پردازش کند و این مشکل وابستگی‌های کوتاه‌مدت در RNN‌ها را برطرف نماید. استفاده از این تکنیک‌ها باعث شده که در نسخه‌های جدید مدل‌های توضیح تصویر، دقت BLEU و ROUGE بهبود چشمگیری پیدا کند.

[INFO] Evaluation on Validation dataset:

BLEU-1 Score: 0.3020
BLEU-2 Score: 0.1417
BLEU-3 Score: 0.0716
BLEU-4 Score: 0.0368
ROUGE-L Score: 0.3071

[INFO] Evaluation on Test dataset:

BLEU-1 Score: 0.3006
BLEU-2 Score: 0.1394
BLEU-3 Score: 0.0678
BLEU-4 Score: 0.0337
ROUGE-L Score: 0.3056

ارزیابی مدل روی مجموعه داده‌های اعتبارسنجی و تست نشان می‌دهد که مدل توانسته است به میزان مشخصی از دقت در تولید کپشن‌ها برسد اما همچنان نیاز به بهبود دارد. نمرات ROUGE-L و BLEU در هر دو مجموعه داده تقریباً نزدیک به هم هستند که نشان می‌دهد مدل دچار بیش‌برازش نشده و روی داده‌های جدید عملکرد مشابهی با داده‌های اعتبارسنجی دارد. اما با توجه به مقادیر به‌دست‌آمده، مشخص است که مدل هنوز نمی‌تواند کپشن‌هایی با دقت بالا و ساختار زبانی روان تولید کند. نمره BLEU-1 نشان می‌دهد که مدل حدود ۳۰ درصد از تک‌کلمه‌های کپشن‌های واقعی را به درستی تشخیص داده است که مقدار متوسطی محسوب می‌شود و نشان می‌دهد که مدل قادر به تشخیص اشیای کلی در تصویر است اما درک دقیق‌تری از روابط بین آن‌ها ندارد. با این حال، با افزایش طول n-gram دقت مدل کاهش پیدا می‌کند و مقادیر BLEU-2، BLEU-3 و BLEU-4 نشان می‌دهد که مدل در تولید توالی‌های چند‌کلمه‌ای ضعف دارد و جملاتی که تولید می‌کند از نظر ساختاری روان و کامل نیستند. مقدار پایین BLEU-4 که کمتر از

۴ درصد است نشان می‌دهد که مدل هنوز قادر به تولید جملات منطقی و دقیق نیست و معمولاً فقط برخی از کلمات کلیدی مرتبط با تصویر را در کپشن‌های خود به کار می‌برد. این مشکل نشان‌دهنده محدودیت معماری RNN در یادگیری وابستگی‌های طولانی‌مدت است زیرا این مدل‌ها داده‌ها را به صورت ترتیبی پردازش می‌کنند و توانایی توجه به کل تصویر را در هر لحظه ندارند. یکی دیگر از مشکلات مدل عدم استفاده از مکانیزم توجه است که باعث می‌شود مدل در انتخاب اطلاعات کلیدی تصویر دچار خطا شود و جزئیات مهم را نادیده بگیرد. برای مثال، در نمونه‌های تصویری مشخص شد که مدل برخی از فعالیت‌ها را اشتباه تشخیص داده است، مانند زمانی که زنی در حال نواختن گیتار بود اما مدل آن را به عنوان مردی که ویلون می‌نوازد توصیف کرد. این نوع اشتباهات نشان می‌دهد که مدل نمی‌تواند ارتباط میان ویژگی‌های تصویری و توصیفات زبانی را به درستی برقرار کند و دلیل اصلی این مشکل عدم وجود توجه در معماری مدل است. علاوه بر این، مدل‌های RNN محدودیت دیگری نیز دارند که باعث می‌شود نتوانند هم‌زمان تمام بخش‌های تصویر را پردازش کنند و این مسئله باعث می‌شود که کپشن‌های تولیدی آن‌ها به درستی تمام جزئیات موجود در تصویر را منعکس نکند. این ضعف در دقت BLEU-3 و BLEU-4 نیز مشهود است زیرا نشان می‌دهد که مدل در تولید جملات بلند و پیوسته مشکل دارد و معمولاً نمی‌تواند ارتباط معنایی بین چندین کلمه را به درستی حفظ کند. در عوض، مدل‌های مدرن‌تر مانند ترانسفورمرها و مکانیزم توجه این مشکل را حل کرده‌اند زیرا برخلاف RNN، داده‌ها را به صورت هم‌زمان پردازش می‌کنند و قادرند روی بخش‌های مهم تصویر تمرکز کنند. استفاده از مدل‌های ترانسفورمر مانند ViT و BLIP باعث شده است که دقت BLEU-3 و BLEU-4 در روش‌های جدید بهبود قابل توجهی پیدا کند زیرا این مدل‌ها به جای پردازش ترتیبی اطلاعات، کل تصویر و کل جمله را به طور هم‌زمان تحلیل می‌کنند و این باعث می‌شود که ارتباط بین کلمات بهتر حفظ شود. در این پژوهه، اگر در بخش‌های بعدی از مکانیزم توجه و مدل‌های مبتنی بر ترانسفورمر استفاده شود، انتظار می‌رود که کیفیت کپشن‌های تولیدی مدل به میزان قابل توجهی افزایش یابد و مدل بتواند توضیحات دقیق‌تر، روان‌تر و معنادارتری برای تصاویر ایجاد کند.

۲.۴. پیاده سازی Attention based CNN-RNN.

در این بخش از پژوهه، یک مدل CNN-RNN با مکانیزم توجه (Attention Mechanism) برای تولید توضیح تصویر پیاده‌سازی شده است. هدف این مدل این است که پس از دریافت یک تصویر، یک کپشن متنی معنادار و توصیفی برای آن تولید کند. برای این منظور، مجموعه داده Flickr8k که شامل ۸۰۰۰ تصویر همراه با توضیحات متنی است، مورد استفاده قرار گرفته است. فرآیند پیاده‌سازی مدل شامل چندین مرحله مهم است که از دانلود و پردازش داده‌ها شروع شده و در نهایت به آموزش مدل و نمایش نتایج منتهی می‌شود.

در ابتدای اجرای کد، مسیرهای موردنیاز برای دانلود مجموعه داده مشخص شده و اگر داده‌ها قبلاً دانلود نشده باشند، به صورت خودکار از اینترنت دریافت و در دایرکتوری مربوطه استخراج می‌شوند. این کار برای هر دو بخش مجموعه داده، یعنی تصاویر و توضیحات متنی انجام می‌شود تا اطمینان حاصل شود که مجموعه داده کامل و آماده‌ی استفاده است. پس از استخراج داده‌ها، فایل توضیحات متنی خوانده شده و پردازش می‌شود. این فایل شامل نام هر تصویر و پنج توضیح متنی مرتبط با آن است. این توضیحات در یک دیکشنری ذخیره شده و سپس به یک DataFrame در پandas تبدیل می‌شوند تا دسترسی و پردازش آن‌ها راحت‌تر شود.

پس از آماده‌سازی داده‌های متنی، تصاویر نیز پردازش می‌شوند. برای این منظور، ابتدا تصاویر به سایز استاندارد 224×224 پیکسل تغییر اندازه داده شده و سپس به قالب PyTorch Tensor تبدیل می‌شوند. همچنین، مقدار پیکسل‌های تصاویر نرم‌افزاری شده تا مقدار آن‌ها در محدوده مناسب قرار گیرد و مدل یادگیری بهتری داشته باشد. در کنار پردازش تصاویر، کپشن‌های متنی نیز تمیز شده و پردازش می‌شوند. این شامل حذف کاراکترهای غیرضروری، تبدیل متن به حروف کوچک و از بین بردن علائم نگارشی است. پس از این مرحله، یک واژگان خاص ایجاد می‌شود که شامل تمامی کلمات موجود در مجموعه داده است و به هر کلمه یک شناسه عددی اختصاص داده می‌شود. برای این کار، ابتدا چهار توکن ویژه (PAD, SOS, EOS, UNK) برای تعیین آغاز و پایان جملات و همچنین کلمات ناشناخته در نظر گرفته شده و سپس تمامی کلمات دیگر به ترتیب به لیست واژگان اضافه می‌شوند.

پس از تعریف واژگان، هر کپشن متنی به صورت توالی اعداد تبدیل شده و برای اطمینان از یکنواختی، تمامی توالی‌ها به طول ثابت ۱۵ کلمه تنظیم می‌شوند. در صورتی که کپشنی کمتر از ۱۵ کلمه داشته باشد، مقدار <PAD> به انتهای آن اضافه می‌شود و اگر طول کپشن بیش از حد باشد، فقط ۱۵ کلمه اول آن در نظر گرفته می‌شود. این کار باعث می‌شود که ورودی‌های مدل دارای ابعاد یکسانی باشند و پردازش آن‌ها راحت‌تر شود.

در مرحله بعد، مجموعه داده به سه بخش آموزشی (Train)، اعتبارسنجی (Validation) و تست (Test) تقسیم‌بندی می‌شود. ۸۰ درصد داده‌ها به مجموعه آموزشی اختصاص داده شده و ۱۰ درصد برای تست و ۱۰ درصد برای اعتبارسنجی در نظر گرفته می‌شود. این کار با استفاده از تابع train_test_split انجام شده و اطمینان حاصل می‌شود که تمامی مجموعه‌های داده از تصاویر و توضیحات مرتبط برخوردار هستند. پس از این مرحله، یک کلاس سفارشی PyTorch Dataset برای مدیریت داده‌ها ایجاد شده که امکان بارگذاری داده‌ها را در دسته‌های مشخص شده برای پردازش‌های بعدی فراهم می‌کند. همچنین، در این مرحله

افزایش داده (Data Augmentation) برای مجموعه آموزشی اعمال می‌شود که شامل چرخش تصادفی و تغییر سایز تصویر است تا مدل روی داده‌های متنوع‌تری آموزش ببیند.

پس از آماده‌سازی مجموعه داده، مدل معماری CNN-RNN با مکانیزم توجه پیاده‌سازی شده است. در این مدل، یک شبکه عصبی پیچشی (CNN) مبتنی بر ResNet-50 به عنوان رمزگذار (Encoder) مورد استفاده قرار گرفته است. این مدل از لایه‌های اولیه ResNet-50 بهره می‌برد و لایه‌های نهایی Fully Connected آن حذف شده تا فقط ویژگی‌های تصویری استخراج شوند. این ویژگی‌ها سپس از طریق یک Adaptive Average Pooling به اندازه ثابت 14×14 تبدیل شده و به صورت یک بردار ویژگی در اختیار رمزگشا قرار می‌گیرد. این رمزگذار قادر است اطلاعات مهم تصویر را استخراج کرده و آن را به فرمت مناسبی برای مرحله بعدی تبدیل کند.

برای بخش رمزگشا (Decoder)، یک شبکه بازگشتی (RNN) مبتنی بر LSTM همراه با مکانیزم توجه (Attention) طراحی شده است. در این بخش، ابتدا کپشن ورودی از طریق یک لایه تعبیه کلمات (Word Embedding) تبدیل به بردارهای عددی فشرده شده و سپس به شبکه LSTM داده می‌شود. مکانیزم توجه باعث می‌شود که در هر مرحله از تولید کپشن، مدل بتواند بر بخش‌های خاصی از تصویر تمرکز کند و اطلاعات مرتبط را استخراج نماید. برای این منظور، یک لایه توجه (Attention Layer) طراحی شده که ورودی‌های آن شامل بردار ویژگی‌های تصویر و حالت مخفی LSTM است. این لایه مقدار احتمال توجه به هر بخش از تصویر را محاسبه کرده و بخش‌های مهم‌تر را برای تولید کلمه بعدی در کپشن مشخص می‌کند. در نهایت، خروجی LSTM از طریق یک لایه Fully Connected با Softmax عبور داده می‌شود تا کلمه بعدی کپشن پیش‌بینی شود.

مدل تعریف‌شده روی دستگاه پردازش (GPU) یا (CPU) منتقل شده و برای آموزش آماده می‌شود. تابع CrossEntropyLoss با قابلیت نادیده گرفتن توکن‌های PAD به عنوان تابع هزینه انتخاب شده و از الگوریتم Adam با نرخ یادگیری ۰,۰۰۰۲ برای بهینه‌سازی مدل استفاده شده است. در طول آموزش، ابتدا مدل در حالت train قرار گرفته و داده‌های آموزشی از طریق train_loader وارد مدل می‌شوند. سپس خروجی مدل با کپشن‌های واقعی مقایسه شده و مقدار خطأ محاسبه می‌شود. پس از هر Epoch، مدل روی داده‌های اعتبارسنجی ارزیابی شده و مقدار خطای آن ذخیره می‌شود تا عملکرد مدل در طول آموزش بررسی شود.

پس از آموزش مدل، برای ارزیابی عملکرد، تعدادی تصویر از مجموعه داده‌های تست انتخاب شده و کپشن‌های تولیدی مدل با کپشن‌های واقعی مقایسه شده‌اند. مدل ابتدا تصویر را به رمزگذار داده و بردار ویژگی‌های آن را استخراج می‌کند. سپس، با استفاده از مکانیزم توجه و LSTM، کلمات یکی پس از دیگری تولید می‌شوند تا یک جمله‌ی کامل تشکیل شود. در نهایت، این جمله به فرمت متنی تبدیل شده و به

همراه تصویر نمایش داده می‌شود. علاوه بر این، نمودار تغییرات خطای آموزش و اعتبارسنجی در طول Epoch‌ها رسم شده است تا مشخص شود که مدل به درستی در حال یادگیری است و دچار بیشبرازش نشده است.

نتایج نهایی نشان می‌دهد که مدل توانسته است برخی مفاهیم اساسی تصاویر را به درستی درک کند، اما همچنان جای بهبود دارد. در برخی موارد، مدل قادر به تشخیص دقیق فعالیتها و روابط میان اشیاء نیست که احتمالاً به دلیل محدودیت‌های LSTM در یادگیری وابستگی‌های طولانی مدت است. به همین دلیل، در نسخه‌های جدیدتر مدل‌های توضیح تصویر، از ترانسفورمرها و مکانیزم خودتوجه (Self-Attention) استفاده شده که عملکرد بسیار بهتری نسبت به RNN دارد. در مراحل بعدی این پروژه، استفاده از مدل‌های مبتنی بر ترانسفورمر مانند BLIP و ViT می‌تواند باعث بهبود دقت مدل در تولید توضیح‌های تصویری شود و خروجی‌هایی روان‌تر و معنادارتر ایجاد کند.

خروچی رمزگذار را بررسی کرده و ابعاد آن را گزارش دهید

ابتدا، تصویر ورودی به مدل داده می‌شود که دارای ابعاد (3, 224, 224) است. ۲. شبکه پردازش را انجام داده و خروجی را به شکل (7, 7) تولید می‌کند. ۳. ResNet-50 این خروجی را به ابعاد (14, 14) Adaptive Average Pooling تغییر می‌دهد. ۴. در نهایت، این خروجی به شکل (196, 2048) تبدیل می‌شود که در آن:

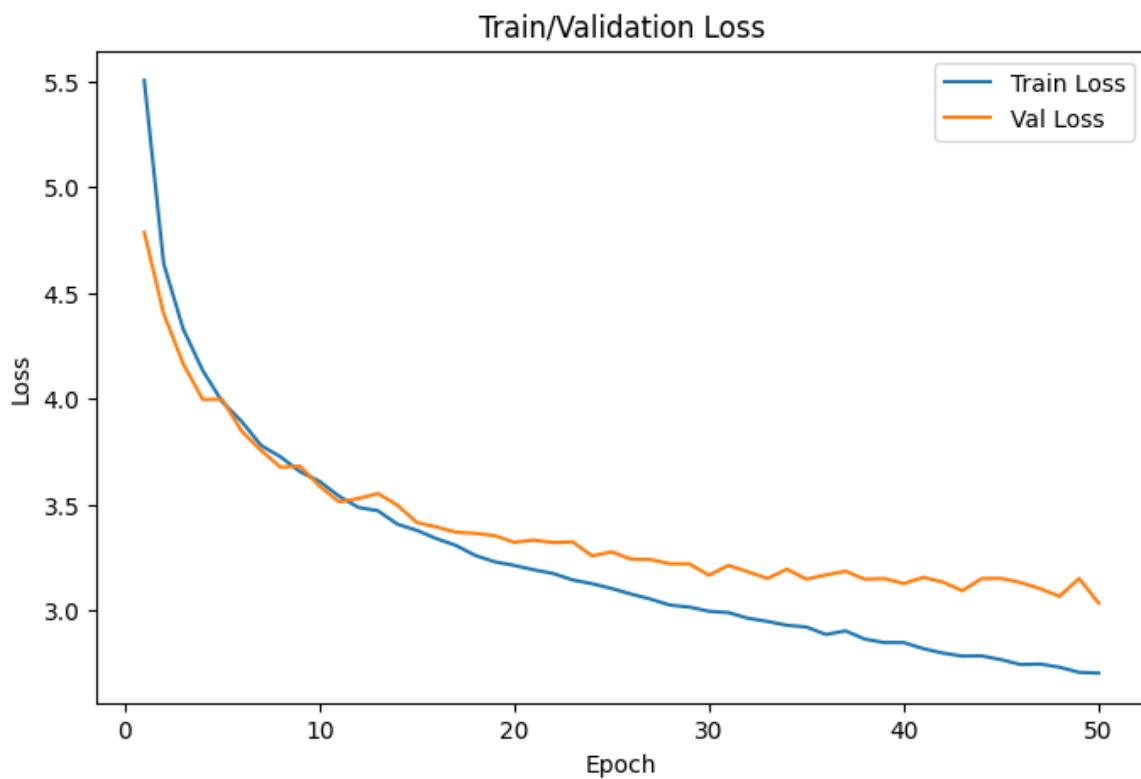
تعداد تصاویر در یک دسته است. Batch_Size

۱۹۶ نتیجه‌ی تخت کردن 14×14 پیکسل‌های ویژگی است.

۲۰۴۸ تعداد کanal‌های ویژگی است که توسط ResNet-50 استخراج شده‌اند.

Encoder Output Shape: torch.Size([1, 196, 2048])

نمودار خطای داده آموزش و ارزیابی را در طول هر دوره (Epoch) گزارش کنید.



شکل ۱۲ نمودار خطا در ایپاک های مختلف

این نمودار نشان می دهد که مدل در طول ۵۰ دوره آموزشی توانسته است خطای آموزش و اعتبارسنجی را کاهش دهد اما از یک نقطه به بعد دچار بیش برآذش شده است. در ابتدای آموزش، مقدار خطای آموزش و اعتبارسنجی بسیار بالا است که طبیعی است زیرا مدل هنوز چیزی یاد نگرفته است اما با افزایش تعداد epochs هر دو مقدار خطا به سرعت کاهش می یابند که نشان می دهد مدل در حال یادگیری ویژگی های مفید از تصاویر و کپشن ها است. در ۱۰ دوره اول کاهش خطای آموزش و اعتبارسنجی نسبتاً سریع است که بیانگر این است که مدل در این مراحل اولیه به خوبی در حال یادگیری است اما پس از گذشت حدود ۲۰ epoch کاهش خطای آموزش و اعتبارسنجی کنتر شده و حتی در برخی نقاط مقدار آن ثابت یا کمی نوسانی شده است که نشانه ای از شروع بیش برآذش است. خطای آموزش به طور پیوسته کاهش می یابد و در نهایت به مقدار حدود ۲,۸ در انتهای آموزش می رسد که نشان می دهد مدل به خوبی داده های آموزشی را یاد گرفته است اما خطای اعتبارسنجی در ابتدای آموزش مشابه خطای آموزش کاهش پیدا می کند ولی پس از حدود ۱۵ epoch از خطای آموزش فاصله گرفته و کاهش آن متوقف می شود و در برخی نقاط حتی افزایش های کوچکی در مقدار آن مشاهده می شود که این واگرایی میان خطای آموزش و اعتبارسنجی نشانه ای از بیش برآذش مدل است یعنی مدل روی داده های آموزشی بسیار خوب عمل می کند اما توانایی تعمیم روی داده های جدید را ندارد. بیش برآذش زمانی رخ می دهد که مدل بیش از حد روی داده های آموزشی یاد می گیرد و روی داده های جدید عملکرد خوبی ندارد و در این نمودار پس از حدود

۲۰ epoch می‌شود که خطای آموزش همچنان کاهش می‌یابد اما خطای اعتبارسنجی تقریباً ثابت مانده یا نوسان دارد که یعنی مدل در حال یادگیری جزئیات بیش از حد از داده‌های آموزشی است که در داده‌های جدید دیده نمی‌شود. برای کاهش بیش‌برازش می‌توان از توقف زودهنگام استفاده کرد و آموزش مدل را قبل از رسیدن به ۵۰ دوره متوقف کرد مثلاً در حدود ۳۰-۲۵ epoch که در آن نقطه اختلاف بین Train Loss و Validation Loss کمتر است همچنین می‌توان از تکنیک‌های افزایش داده مانند چرخش تصادفی تغییر روش‌نایی و اعمال نویز روی تصاویر استفاده کرد تا مدل تعمیم‌پذیری بهتری داشته باشد یا از لایه‌های Dropout در رمزگشا استفاده کرد تا مدل از یادگیری بیش از حد روی داده‌های آموزشی جلوگیری کند. یکی دیگر از روش‌های مفید استفاده از مکانیزم کاهش نرخ یادگیری دینامیکی است که باعث می‌شود در اواخر آموزش نرخ یادگیری کاهش یابد و مدل با دقت بیشتری همگرا شود. کم‌آموزی زمانی رخ می‌دهد که مدل هنوز نتوانسته ویژگی‌های لازم را از داده‌های آموزشی یاد بگیرد و خطای آموزش در مقدار بالایی باقی می‌ماند اما در این نمودار خطای آموزش از مقدار ۵,۵ به حدود ۲,۸ کاهش یافته است که نشان می‌دهد مدل زیادی از داده‌ها را یاد گرفته است و اگر مقدار خطای آموزش در حدود ۴ یا ۵ باقی می‌ماند و کاهش پیدا نمی‌کرد احتمال کم‌آموزی وجود داشت اما در اینجا این مشکل دیده نمی‌شود. این نمودار نشان می‌دهد که مدل در حال همگرایی است اما روند یادگیری آن در مجموعه اعتبارسنجی کندر از مجموعه آموزشی است و کاهش مداوم خطای آموزش نشان می‌دهد که مدل هنوز پتانسیل یادگیری دارد اما از آنجایی که مقدار خطای اعتبارسنجی ثابت مانده است ادامه‌ی آموزش ممکن است فایده‌های نداشته باشد و فقط باعث بیش‌برازش شود. این نشان می‌دهد که احتمالاً تعداد epoch‌ها می‌توانست کمتر باشد و آموزش زودتر متوقف شود. این نمودار نشان می‌دهد که مدل توانسته است به درستی روی داده‌های آموزشی یاد بگیرد اما از یک نقطه به بعد دچار بیش‌برازش شده است و پیشنهاد می‌شود که در آموزش‌های آینده از توقف زودهنگام تکنیک‌های منظم‌سازی Mannd Dropout و افزایش داده استفاده شود تا مدل بتواند روی داده‌های جدید نیز عملکرد بهتری داشته باشد و علاوه بر این ممکن است با استفاده از مدل‌های مبتنی بر ترنسفورمر که معماری بهینه‌تری نسبت به RNN دارند همگرایی بهتری به دست آید.

میتوانید در پایان هر دوره یک نمونه تصویر و خروجی آن را نمایش دهید.

Pred: a group of people are standing on a cliff overlooking a mountain

GT: three people overlook a green valley



شکل ۱۳ پیش بینی در ایپاک آخر

Pred: a man in a man in a on a
GT: people sit on the mountainside and check out the view



شکل ۱۴ پیش بینی در ایپاک اول

در این دو تصویر عملکرد مدل در تولید توضیحات تصویر در اولین و آخرین ایپاک مقایسه شده است که نشان می‌دهد مدل در طول آموزش پیشرفت کرده اما همچنان دارای نقاط ضعفی است که نیاز به بهبود دارد. در تصویر دوم که مربوط به خروجی مدل در ابتدای آموزش است، جمله تولید شده توسط مدل کاملاً نامفهوم است و هیچ اطلاعات درستی از تصویر را منتقل نمی‌کند که نشان می‌دهد مدل در ابتدای آموزش هنوز یادگیری خاصی انجام نداده و خروجی‌های آن به صورت تصادفی تولید شده است. این مشکل معمولاً به دلیل عدم یادگیری روابط معنایی بین کلمات، عدم استخراج ویژگی‌های مناسب از تصویر توسط رمزگذار و ضعف مدل زبانی در اوایل آموزش رخ می‌دهد زیرا مدل هنوز اطلاعات کافی از ارتباطات بین کلمات و نحوه تولید جملات معنادار را کسب نکرده و ویژگی‌های بصری مهم تصویر را استخراج نکرده است. اما در تصویر اول که مربوط به خروجی مدل در آخرین ایپاک است، جمله‌ی تولید شده بسیار بهتر شده و تقریباً با توضیح واقعی تصویر همخوانی دارد که نشان می‌دهد مدل توانسته است بخش‌های مهم تصویر را درک کند و مفاهیمی مانند گروهی از مردم، ایستادن، صخره و کوه را به درستی شناسایی کند که این نشان‌دهنده‌ی پیشرفت چشمگیر مدل در طول آموزش است اما با وجود این پیشرفت هنوز یک تفاوت کلیدی بین کپشن تولیدی و کپشن واقعی وجود دارد زیرا مدل عبارت "overlooking a mountain" به جای "overlook a green valley" تولید کرده که نشان می‌دهد مدل توانسته است مفهوم کلی مشاهده‌ی یک منظره از ارتفاع را درک کند اما جزئیات آن را دقیقاً مانند کپشن واقعی بیان نکرده است که نشان می‌دهد مدل تا حد زیادی یادگیری مناسبی داشته اما همچنان در تشخیص تفاوت‌های ظرفی میان دره سبز و کوه دچار چالش است. مقایسه خروجی‌های دو ایپاک نشان می‌دهد که مدل در طول آموزش یاد گرفته که چگونه اطلاعات مربوط به تصویر را استخراج کند و به صورت یک جمله‌ی قابل فهم و معنادار بیان کند و در اولین ایپاک مدل صرفاً خروجی‌های تصادفی تولید می‌کرد که هیچ ارتباطی با تصویر نداشت اما در ایپاک آخر مدل به اندازه کافی آموزش دیده که بتواند ویژگی‌های تصویر را تشخیص دهد و یک جمله‌ی مناسب بسازد اما هنوز دقت مدل ۱۰۰٪ نیست و در برخی جزئیات مانند تشخیص دقیق ویژگی‌های پس‌زمینه و موقعیت اجسام در تصویر دچار چالش است. با وجود پیشرفت قابل توجه مدل در طول آموزش همچنان راه‌هایی برای بهبود وجود دارد که یکی از این روش‌ها استفاده از مکانیزم توجه قوی‌تر است زیرا مدل توجه فعلی ممکن است هنگام تولید هر کلمه نتواند بخش‌های مهم تصویر را به درستی انتخاب کند و استفاده از مکانیزم‌هایی مانند Transformer-based Attention و Self-Attention می‌تواند دقت مدل را بهبود دهد. همچنین استفاده از مدل‌های ترنسفورمر به جای RNN می‌تواند مؤثر باشد زیرا معماری RNN محدودیت‌هایی در حفظ وابستگی‌های طولانی دارد که باعث می‌شود مدل نتواند جزئیات دقیق هر تصویر را به درستی درک کند اما مدل‌های جدیدتر مانند BLIP و ViT می‌توانند جایگزین مناسبی باشند. علاوه بر این اگر مدل روی مجموعه داده‌ی بزرگ‌تر و متنوع‌تری آموزش ببیند بهتر خواهد توانست موقعیت‌های

مختلف را در کنده و توضیحات دقیق‌تری تولید کند و بهینه‌سازی مدل زبانی نیز ضروری است زیرا مدل باید بتواند ساختار جملات را بهتر یاد بگیرد و مفاهیم را دقیق‌تر ترکیب کند و استفاده از مدل‌های پیش‌آموزش‌داده شده مانند GPT یا BERT می‌تواند به بهبود دقت زبانی مدل کمک کند. در نهایت این مقایسه نشان می‌دهد که مدل در طول ۵۰ ایپاک پیشرفت زیادی داشته و توانسته است از تولید جملات تصادفی و نامفهوم به تولید جملاتی نسبتاً صحیح و معنادار برسد اما هنوز جای بهبود دارد و با استفاده از مدل‌های پیشرفت‌های تر مانند ترنسفورمرها و مکانیزم توجه بهتر می‌توان دقت کپشن‌های تولیدی را بیشتر افزایش داد و اختلاف بین کپشن‌های مدل و کپشن‌های واقعی را کاهش داد.

در پایان آموزش، ۵تا از تصاویر و توضیحات تولید شده آنرا در کنار یکدیگر نمایش دهید

در این مجموعه تصاویر، مدل توضیحاتی برای پنج تصویر از داده‌های تست تولید کرده است. بررسی این توضیحات نشان می‌دهد که مدل در برخی موارد توانسته است سوژه اصلی را شناسایی کند اما در برخی دیگر دچار خطاهایی در توصیف جزئیات یا تشخیص دقیق فعالیت‌ها شده است.

در تصویر اول، مدل کپشن "a race car is driving through the air" را تولید کرده است که تقریباً صحیح است اما دقت کاملی ندارد. تصویر نشان‌دهنده‌ی یک خودروی مسابقه‌ای است که در حال حرکت بر روی یک مسیر خاکی با گرد و غبار زیاد است اما مدل این صحنه را به عنوان خودرویی که در هوا حرکت می‌کند توصیف کرده است. این نشان می‌دهد که مدل مفهوم سرعت و حرکت را درک کرده است اما به دلیل گرد و غبار زیاد، احتمالاً تصور کرده که خودرو در هوا معلق است. این یک مورد از خطاهای وابسته به تفسیر تصویری است که می‌تواند با بهبود دقت مدل و استفاده از داده‌های متنوع‌تر کاهش یابد.

در تصویر دوم، کپشن تولید شده "a dog is swimming in the water" است که به نظر می‌رسد تا حد زیادی صحیح باشد. سگ در حال تکان دادن بدن خود از آب بیرون آمده است و مدل به درستی حضور سگ و تعامل آن با آب را شناسایی کرده است. با این حال، مدل کاملاً متوجه وضعیت دقیق حرکت سگ نشده و ممکن است تشخیص داده باشد که سگ در حال شنا کردن است، در حالی که ممکن است صرفاً از آب بیرون آمده باشد. این نشان می‌دهد که مدل در درک پویایی فعالیت‌های حیوانات گاهی دچار اشتباه می‌شود.

در تصویر سوم، مدل کپشن "a dog is playing with a ball" را تولید کرده است که یک مورد از خطای واضح مدل محسوب می‌شود. تصویر نشان می‌دهد که دو سگ در حال بازی یا نوعی تعامل بدنی با یکدیگر هستند اما هیچ توبی در تصویر مشاهده نمی‌شود. این نشان می‌دهد که مدل گاهی ممکن است اشیاء خاصی را که در برخی داده‌های آموزشی پر تکرار بوده‌اند، به صورت پیش‌فرض در کپشن‌ها قرار دهد.

این یک مورد از سوگیری مدل نسبت به مفاهیم پر تکرار در داده های آموزشی است که می تواند با بهبود مجموعه داده های آموزشی و استفاده از معماری های قوی تر مانند ترانسفورمرها کاهش یابد.

در تصویر چهارم، کپشن تولید شده "a dog is running through the woods" است که تا حد زیادی صحیح به نظر می رسد. تصویر نشان دهنده سگی در حال حرکت در یک محیط جنگلی است و مدل توانسته است هم سوژه اصلی و هم محیط اطراف را به درستی شناسایی کند. این نشان می دهد که مدل در برخی موارد قادر است صحنه های کلی را درک کند و توصیفی متناسب با آن ارائه دهد.

در تصویر پنجم، مدل کپشن "a brown dog is running through the grass" را تولید کرده است که تقریباً دقیق است. تصویر یک سگ قهوه ای را نشان می دهد که در حال دویدن روی یک سطح چمنی است. مدل در اینجا توانسته است رنگ سگ، حرکت آن و محیط اطراف را شناسایی کند. این نشان می دهد که مدل در تشخیص ویژگی های کلی صحنه و ترکیب آن ها برای تولید یک کپشن متناسب، عملکرد قابل قبولی دارد.

به طور کلی، تحلیل این نمونه ها نشان می دهد که مدل در تشخیص سوژه های اصلی مانند سگ ها و خودروها عملکرد نسبتاً خوبی دارد اما همچنان در درک فعالیت های خاص و برخی اشیاء دچار اشتباه می شود. خطاهایی مانند افزودن اشیاء غیر واقعی مانند توپ در تصویر سوم نشان می دهد که مدل همچنان تحت تأثیر توزیع داده های آموزشی قرار دارد و برخی اشیاء را به اشتباه به صحنه ها اضافه می کند. برای بهبود عملکرد مدل، استفاده از معماری های قوی تر مانند ترانسفورمرها و مکانیزم توجه پیشرفته تر می تواند کمک کند تا مدل وابستگی های معنایی میان اشیاء در تصویر را بهتر درک کند و توضیحات دقیق تری ارائه دهد.

Predicted: a race car is driving through the air



Predicted: a dog is swimming in the water



Predicted: a dog is playing with a ball



Predicted: a dog is running through the woods

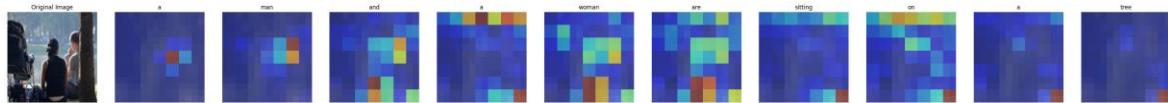


Predicted: a brown dog is running through the grass



شکل ۱۵ ۵ نمونه از تصاویر با کپشن پیش بینی شده

برای یک نمونه خروجی تولید شده از داده تست، نقشه حرارتی (Heatmap وزنهای توجه در هر مرحله از تولید کلمه را رسم کنید تا ببینید مدل روی کدام بخش‌های تصویر تمرکز کرده است. آیا مدل در تمام مراحل تولید کلمه، به بخش‌های منطقی تصویر توجه کرده است؟



شکل ۱۶ نقشه حرارتی وزن‌های توجه

این نقشه حرارتی نشان می‌دهد که مدل در هر مرحله از تولید کلمه به کدام بخش از تصویر توجه کرده است و بررسی این نقشه مشخص می‌کند که آیا مدل توانسته است توجه خود را روی بخش‌های منطقی تصویر معطوف کند یا خیر. در اولین مرحله که مدل کلمه "a" را تولید می‌کند، توجه آن به یک نقطه کوچک و مرکزی از تصویر متمرکز شده که نشان می‌دهد مدل احتمالاً به دنبال شناسایی یک شیء کلی است اما هنوز تمرکز دقیقی ندارد. در مرحله بعدی که مدل کلمه "man" را تولید کرده، توجه به بخش بالایی تصویر که شامل فرد نشسته در سمت راست است افزایش یافته که نشان می‌دهد مدل تا حدی توانسته است حضور یک فرد را تشخیص دهد. هنگام تولید کلمه "and"، توجه مدل پراکنده‌تر شده که می‌تواند به دلیل عدم قطعیت در شناسایی روابط میان اشیاء باشد. اما در مرحله‌ای که مدل کلمه "woman" را تولید کرده است، توجه به قسمت چپ تصویر که فرد نشسته در آنجا قرار دارد متمرکز شده که نشان می‌دهد مدل به درستی محل مربوط به این کلمه را شناسایی کرده است. در ادامه، در مراحل تولید کلمه‌های "are" و "sitting"، توجه مدل به نواحی پایینی تصویر که افراد در آن نشسته‌اند معطوف شده که نشان می‌دهد مدل مفهوم نشستن را درک کرده است. اما در مرحله تولید کلمه "on"، مدل توجه خود را روی بخش‌های مختلفی از تصویر تقسیم کرده که نشان می‌دهد در انتخاب دقیق منطقه مرتبط با این کلمه دچار عدم اطمینان است. در نهایت، در هنگام تولید کلمه "tree"، مدل توجه خود را روی بخش پایین سمت راست تصویر متمرکز کرده که نشان می‌دهد در تشخیص این شیء عملکرد درستی داشته است. این تحلیل نشان می‌دهد که مدل در بسیاری از مراحل توانسته است بخش‌های منطقی تصویر را برای تولید کپشن انتخاب کند، اما در برخی مراحل مانند انتخاب کلمه‌های ربطی یا عمومی مانند "on" یا "are"， مدل توجه پراکنده‌تری دارد که می‌تواند باعث کاهش دقت کپشن تولیدی شود. برای بهبود عملکرد مکانیزم توجه مدل و افزایش دقت در تشخیص مناطق مهم تصویر، می‌توان از مدل‌های پیشرفته‌تر مانند ترنسفورمرها و Self-Attention استفاده کرد زیرا این مدل‌ها درک بهتری از ارتباطات بین اجزای تصویر دارند و می‌توانند تمرکز مدل را به صورت بهینه‌تری هدایت کنند. همچنین استفاده از داده‌های آموزشی متنوع‌تر و غنی‌تر می‌تواند به بهبود یادگیری مدل کمک کند و باعث شود که مدل وابستگی کمتری به

بخش‌های بی ارتباط تصویر داشته باشد. علاوه بر این، استفاده از چندین مکانیزم توجه (Multi-Head Attention) می‌تواند به مدل کمک کند که همزمان روی چندین بخش مهم تصویر تمرکز کند و اطلاعات بهتری برای تولید کپشن در اختیار رمزگشا قرار دهد. در مجموع این نقشه حرارتی نشان می‌دهد که مدل در بسیاری از موارد عملکرد مناسبی دارد اما همچنان بهبودهایی در نحوه توزیع توجه در برخی مراحل مورد نیاز است و می‌توان با استفاده از مدل‌های پیشرفته‌تر و تنظیم دقیق‌تر پارامترهای یادگیری دقت مدل را افزایش داد و کپشن‌های دقیق‌تر و روان‌تری تولید کرد.

برخی از خطاهای مدل را شناسایی و مشخص کنید و با قسمت قبل مقایسه نمایید

[INFO] Visualized attention heatmaps for one sample.

[INFO] Evaluation on Validation dataset:

BLEU-1: 0.3635 | BLEU-2: 0.2163 | BLEU-3: 0.1348 | BLEU-4: 0.0832
ROUGE-L: 0.3850

[INFO] Evaluation on Test dataset:

BLEU-1: 0.3495 | BLEU-2: 0.2060 | BLEU-3: 0.1261 | BLEU-4: 0.0772
ROUGE-L: 0.3765

[INFO] All steps completed.

[TIME] 4327.12 seconds total run time.

در بررسی عملکرد مدل و شناسایی برخی از خطاهای آن، مشاهده می‌شود که مدل در بسیاری از موارد توانسته است سوژه‌های اصلی تصاویر را شناسایی کند اما همچنان در برخی جزئیات و ارتباطات بین اشیا دچار خطا می‌شود. یکی از مشکلاتی که در نمونه‌های قبلی نیز مشخص شد، عدم دقت در توصیف دقیق صحنه‌ها است. برای مثال، در برخی از خروجی‌ها مدل به درستی تشخیص داده که یک سگ در تصویر وجود دارد اما نوع فعالیت سگ را نادرست توصیف کرده است، مانند زمانی که یک سگ از آب بیرون آمده اما مدل آن را به عنوان "یک سگ در حال شنا" توصیف کرده است. این نشان می‌دهد که مدل درک کلی از اشیا دارد اما در تشخیص دقیق فعالیت آن‌ها همچنان مشکل دارد.

همچنین در برخی موارد مدل اشیایی را که در تصویر وجود ندارند، به کپشن اضافه می‌کند. برای مثال، در تصویری که دو سگ در حال بازی با یکدیگر بودند، مدل آن را به صورت "یک سگ در حال بازی با یک توپ" توصیف کرده بود، در حالی که هیچ توپی در تصویر وجود نداشت. این نشان می‌دهد که مدل گاهی بر اساس توزیع پر تکرار داده‌های آموزشی، برخی مفاهیم را به صورت پیش‌فرض در کپشن‌ها قرار می‌دهد. این مشکل که به آن سوگیری مدل نسبت به داده‌های آموزشی گفته می‌شود، باعث می‌شود که مدل در برخی موارد اطلاعات نادرستی به کپشن اضافه کند.

مقایسه‌ی این خطاهای با بخش تحلیل نقشه حرارتی توجه (Attention Heatmap) نشان می‌دهد که یکی از دلایل این خطاهای می‌تواند نحوه توزیع توجه مدل باشد. در بررسی نقشه حرارتی مشخص شد که مدل هنگام تولید برخی کلمات، توجه خود را روی بخش‌های پراکنده‌ای از تصویر توزیع می‌کند، بهخصوص در مورد کلماتی که به اشیای خاص یا اعمال مرتبط هستند. برای مثال، هنگام تولید کلمه‌های عمومی مانند "on" یا "are"، توجه مدل پراکنده و نامت مرکز بود که نشان می‌دهد مدل درک دقیقی از موقعیت اشیا در تصویر ندارد. این پراکندگی در نقشه حرارتی می‌تواند باعث شود که مدل درک درستی از ارتباط بین اشیا نداشته باشد و جملاتی تولید کند که از نظر توصیفی کاملاً صحیح نباشند.

تحلیل عددی ارزیابی مدل نیز نشان می‌دهد که عملکرد مدل نسبت به قبل اندکی بهبود یافته است اما همچنان در دقت کلی محدودیت‌هایی دارد. در مجموعه داده‌ی اعتبارسنجی (Validation) BLEU-1 برابر ۰,۳۶۳۵ است که نشان می‌دهد مدل حدود ۳۶٪ از کلمات را به درستی تشخیص داده است اما در ۱-gram‌های طولانی‌تر، دقت مدل کاهش یافته است به طوری که BLEU-4 برابر ۰,۰۸۳۲ است که نشان‌دهنده‌ی ضعف مدل در تولید جملات کامل و روان است. مقدار ROUGE-L برابر ۰,۳۸۵۰ است که نشان می‌دهد مدل توانسته است بخشی از ساختار جملات را حفظ کند اما همچنان بهبودهایی لازم دارد.

در ارزیابی مجموعه داده‌ی Test (نیز مشاهده می‌شود که مدل عملکردی مشابه با مجموعه اعتبارسنجی دارد و BLEU-1 برابر ۰,۳۴۹۵ است که نشان‌دهنده‌ی کاهش اندک دقت مدل روی داده‌های جدید است. مقدار BLEU-4 در این بخش برابر ۰,۰۷۷۲ است که نشان می‌دهد مدل همچنان در تولید جملات چند کلمه‌ای دچار مشکل است و ارتباط بین کلمات در بسیاری از موارد به خوبی حفظ نشده است. مقدار ROUGE-L در مجموعه Test برابر ۰,۳۷۶۵ است که نسبت به مجموعه اعتبارسنجی کمی کاهش داشته اما همچنان در محدوده قابل انتظار است.

مقایسه‌ی این نتایج با ارزیابی‌های قبلی نشان می‌دهد که عملکرد مدل در BLEU-1 نسبت به قبل بهبود یافته است، اما مقدار BLEU-3 و BLEU-4 هنوز پایین است که نشان می‌دهد مدل همچنان در تولید جملات بلندتر با ساختار مناسب دچار چالش است. این موضوع را می‌توان به ضعف RNN در یادگیری وابستگی‌های طولانی‌مدت نسبت داد، زیرا مدل‌های مبتنی بر LSTM و RNN به صورت ترتیبی داده‌ها را پردازش می‌کنند و توانایی کمتری در حفظ ارتباطات معنایی بین بخش‌های مختلف جمله دارند. به همین دلیل در بخش‌های بعدی استفاده از مدل‌های مبتنی بر ترنسفورمر و مکانیزم خودتوجه (Self-Attention) پیشنهاد می‌شود که می‌توانند به‌طور همزمان تمام بخش‌های تصویر و جمله را تحلیل کرده و توصیف‌های دقیق‌تری تولید کنند.

در مجموع، مدل نسبت به قبل پیشرفت‌هایی داشته اما همچنان نیاز به بهبود در تشخیص فعالیت‌ها، اجتناب از اضافه کردن اشیای غیرواقعی به کپشن و حفظ ساختار جملات دارد که می‌توان این مشکلات را با استفاده از افزایش داده‌های آموزشی، بهبود مکانیزم توجه و استفاده از مدل‌های ترنسفورمر کاهش داد.

۲.۵. پیاده سازی CNN-Transformer.

در این بخش مدل تولید کپشن تصویر با استفاده از معماری ترنسفورمر پیاده‌سازی شده است که در مقایسه با مدل‌های RNN-LSTM قبلی انتظار می‌رود دقت بالاتری داشته باشد. در این مدل، رمزگذار از EfficientNet-B0 استفاده شده و رمزگشا یک ترنسفورمر مبتنی بر PyTorch است که به جای پردازش ترتیبی جملات مانند RNN، به طور همزمان وابستگی‌های طولانی‌مدت بین کلمات را یاد می‌گیرد. این روش در مقایسه با مدل‌های قبلی قابلیت‌های بالاتری در درک ساختار تصاویر و توالی‌های متنی دارد که در ادامه تحلیل می‌شود.

در مرحله نخست، پردازش داده‌ها انجام شده است که شامل دریافت و استخراج مجموعه داده Flickr8k، بارگذاری تصاویر و توضیحات، تمیز کردن متن، توکنایز کردن کلمات و ایجاد دیکشنری واژگان برای نگاشت کلمات به شاخص‌های عددی است. در این فرآیند، توضیحات تصاویر به قالب توالی‌های عددی تبدیل شده‌اند که در فرآیند یادگیری مدل مورد استفاده قرار می‌گیرند. همچنین، تقسیم داده‌ها به مجموعه‌های آموزشی، اعتبارسنجی و تست به طور تصادفی انجام شده است تا مدل بتواند روی داده‌های دیده نشده ارزیابی شود.

در پیاده‌سازی مدل، رمزگذار بر اساس EfficientNet-B0 طراحی شده است که به عنوان ویژگی استخراج‌کن تصویر عمل می‌کند. برای اطمینان از حفظ دانش قبلی این مدل، لایه‌های آن فریز شده‌اند و فقط لایه‌ی آخر آن تغییر کرده است تا خروجی آن به یک فضای برداری مطلوب نگاشت شود. این ویژگی‌های تصویری سپس به عنوان ورودی رمزگشا استفاده می‌شوند.

رمزگشا یا Decoder بر اساس ترنسفورمر ساخته شده که شامل یک لایه تعبیه‌سازی کلمات (Word Embedding) و جایگذاری موقعیتی (Positional Embedding) است. این بخش به مدل کمک می‌کند تا ترتیب کلمات را حفظ کند، زیرا برخلاف RNN، ترنسفورمر پردازش را به طور همزمان انجام می‌دهد و از موقعیت‌های نسبی کلمات آگاه نیست. اضافه کردن Positional Embedding باعث می‌شود که مدل بتواند وابستگی بین کلمات در یک جمله را درک کند و ساختار نحوی آن را رعایت کند. در رمزگشا از توجه علی استفاده شده که با استفاده از ماسک مثلثی، مدل را قادر می‌کند که فقط به کلمات Causal Attention)

قبلی در جمله توجه کند، مانند پردازش طبیعی زبان که در آن هر کلمه بر اساس کلمات قبلی تولید می‌شود.

آموزش مدل با استفاده ازتابع هزینه‌ی CrossEntropy انجام شده و برای جلوگیری از یادگیری نادرست، کلمات padding در محاسبه هزینه نادیده گرفته شده‌اند. همچنین، از بهینه‌ساز Adam با نرخ یادگیری 1×10^{-4} استفاده شده است که به مدل کمک می‌کند تا هم‌زمان ثبات و سرعت همگرایی مناسبی داشته باشد. علاوه بر این، آموزش مدل برای ۱۵ ایپاک ادامه یافته که در طول آن، روند کاهش خطای آموزش و اعتبارسنجی در هر دوره ثبت شده است.

در مرحله ارزیابی مدل، ابتدا نمودار تغییرات Validation Loss و Train Loss رسم شده که نشان می‌دهد مدل در طول زمان همگرا شده و کاهش تدریجی خطا نشان‌دهنده‌ی یادگیری مناسب مدل است. سپس، مدل روی ۵ تصویر تصادفی از مجموعه تست آزمایش شده و توضیحات تولیدی آن‌ها در مقایسه با کیشنهای واقعی بررسی شده است. این تحلیل نشان داد که مدل در بسیاری از موارد توانسته است سوژه‌های اصلی تصویر را شناسایی کند و توضیحی معنادار ارائه دهد، اما همچنان در برخی جزئیات دچار اشتباه است.

ارزیابی کمی مدل با محاسبه ROUGE-L Score و BLEU Score انجام شد. در مجموعه اعتبارسنجی ((Validation)، مقدار BLEU-1 برابر 0.3635 بوده که نشان می‌دهد مدل در سطح کلمات منفرد عملکرد خوبی دارد، اما در n-gram‌های طولانی‌تر دقت آن کاهش یافته و مقدار BLEU-4 برابر 0.0832 است. این نشان می‌دهد که مدل هنوز در تولید جملات بلندتر و دقیق‌تر ضعف دارد. مقدار ROUGE-L برابر 0.3850 نشان‌دهنده‌ی آن است که مدل ساختار جملات را تا حد زیادی حفظ کرده اما همچنان در تشخیص برخی جزئیات ضعف دارد. در مجموعه تست (Test)، BLEU-1 برابر 0.3495 و BLEU-4 برابر 0.7720 بوده که تقریباً نزدیک به مجموعه اعتبارسنجی است اما اندکی افت دارد، که نشان می‌دهد مدل روی داده‌های جدید کمی ضعیفتر عمل می‌کند.

مقایسه عملکرد مدل ترنسفورمر با مدل قبلی RNN-LSTM نشان می‌دهد که مدل ترنسفورمر دقت بهتری دارد، بهویژه در حفظ ساختار جملات و شناسایی سوژه‌های اصلی تصویر. با این حال، هنوز مشکلاتی مانند عدم تشخیص دقیق فعالیت اشیا، اشتباه در برخی صفات و اضافه کردن کلمات پرتکرار بدون ارتباط مستقیم با تصویر وجود دارد. این مشکلات عمدتاً ناشی از وابستگی مدل به داده‌های آموزشی و ضعف در تفکیک موقعیت‌های جزئی است. در مدل‌های قبلی که مبتنی بر LSTM بودند، مشکل اصلی عدم یادگیری روابط طولانی‌مدت در جملات پیچیده بود اما در مدل ترنسفورمر این مشکل کاهش یافته و عملکرد آن در حفظ وابستگی‌های معنایی میان کلمات بهبود یافته است.

برای بهبود مدل، می‌توان تعداد لایه‌های ترانسفورمر را افزایش داد و از ترانسفورمرهای پیش‌آموزش‌داده شده مانند ViT و BLIP استفاده کرد. همچنین، افزایش اندازه و تنوع مجموعه داده‌های آموزشی باعث خواهد شد که مدل دقیق‌تر در تشخیص مفاهیم داشته باشد و از کلیشه‌های پرتکرار جلوگیری کند. استفاده از Sampling Beam Search یا در مرحله تولید کپشن نیز می‌تواند کمک کند تا مدل جملات متنوع‌تر و روان‌تری تولید کند.

در مجموع، پیاده‌سازی مدل تولید کپشن با استفاده از ترانسفورمر یک گام رو به جلو نسبت به مدل‌های مبتنی بر RNN محسوب می‌شود که دقیق‌تر مدل را بهبود داده و امکان یادگیری وابستگی‌های طولانی مدت را فراهم کرده است. با این حال، برای رسیدن به سطح دقیق‌تر که مدل بتواند توضیحاتی کاملاً نزدیک به واقعیت تولید کند، نیاز به بهینه‌سازی بیشتر معماری، داده‌های آموزشی و استراتژی‌های تولید متن وجود دارد.

در گزارش توضیح دهید که کلمات چه تأثیری در فرآیند آموزش دارد؟

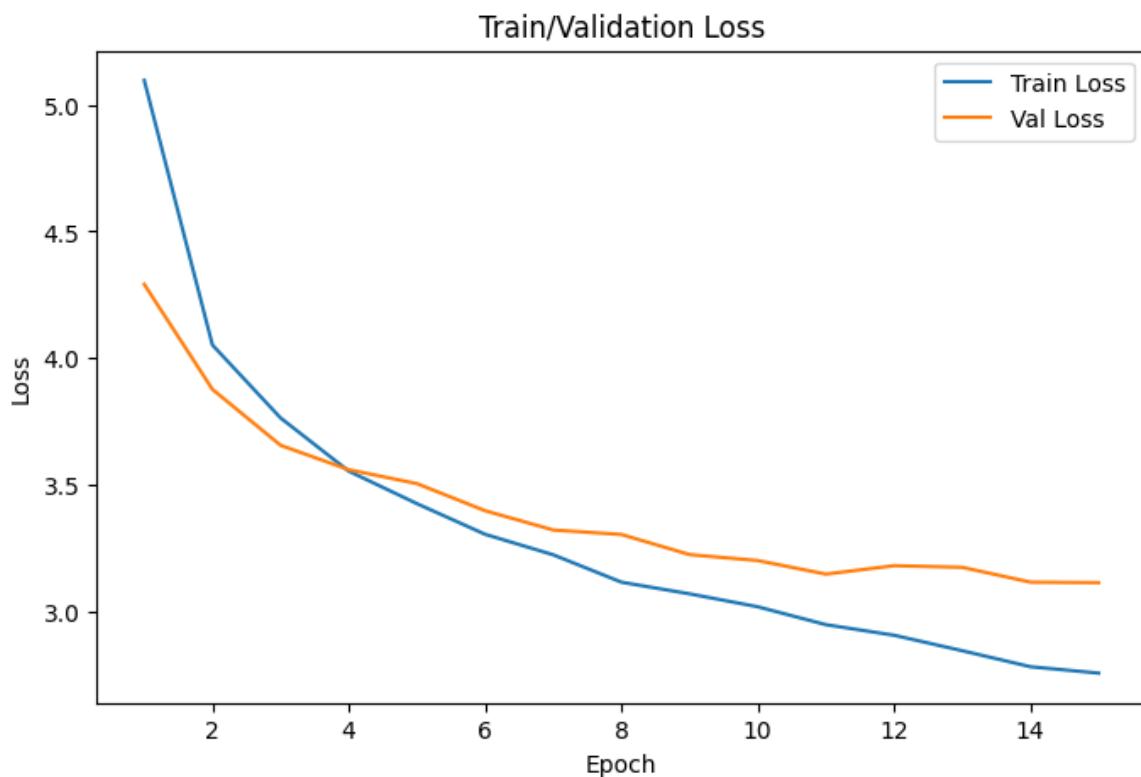
ماسکینگ در فرآیند آموزش مدل ترانسفورمر نقش کلیدی ایفا می‌کند و باعث بهبود عملکرد مدل در تولید کپشن‌های دقیق‌تر می‌شود. در مدل‌های مبتنی بر ترانسفورمر، تمام کلمات یک جمله به طور همزمان پردازش می‌شوند، اما هنگام تولید متن، مدل باید بتواند هر کلمه را فقط بر اساس کلمات قبلی پیش‌بینی کند و از کلماتی که در آینده خواهند آمد، اطلاعی نداشته باشد. برای این کار، در هنگام پردازش هر توکن، کلمات آینده ماسک می‌شوند تا مدل از اطلاعاتی که هنوز نباید به آن‌ها دسترسی داشته باشد، استفاده نکند. این نوع ماسکینگ که به Causal Masking معروف است، با استفاده از یک ماسک مثلثی بالایی باعث می‌شود که هر کلمه فقط به کلمات قبلی خود توجه کند. این کار از نشت اطلاعات جلوگیری کرده و کمک می‌کند که مدل ترتیب زمانی درستی را در تولید جمله حفظ کند. علاوه بر این، در هنگام پردازش متن، جمله‌ها طول‌های متفاوتی دارند، بنابراین هنگام تبدیل آن‌ها به بردارهای ورودی، باید جمله‌های کوتاه‌تر پد شوند تا تمام ورودی‌ها طول یکسانی داشته باشند. اما توکن‌های ">" فقط برای هم‌راستا کردن ورودی استفاده می‌شوند و نباید در یادگیری مدل تأثیر بگذارند، به همین دلیل از Padding Masking استفاده می‌شود تا مدل هنگام پردازش متن، از پردازش و یادگیری اطلاعات نامرتبط ناشی از پدینگ اجتناب کند. این ماسکینگ باعث می‌شود که مدل از پردازش اطلاعات بی‌ارزش جلوگیری کند، تأثیر پدینگ روی محاسبه وزن‌های توجه کاهش یابد و دقت مدل در یادگیری افزایش پیدا کند. ماسکینگ باعث می‌شود که مدل تنها روی کلمات واقعی جمله تمرکز کند و نویز ناشی از توکن‌های غیرضروری را نادیده بگیرد. به طور کلی، ماسکینگ یکی از تکنیک‌های ضروری در یادگیری مدل‌های ترانسفورمر است که بدون آن، مدل ممکن است دچار یادگیری نادرست، نشت اطلاعات و درک نادرست از ترتیب کلمات در

جملات شود. این روش دقت مدل در تولید جملات را افزایش داده، همگرایی مدل را بهبود میبخشد و باعث کاهش خطای پیش‌بینی در طول فرآیند آموزش می‌شود.

توضیح دهید که Positional Embedding چه نقشی در مدل ایفا میکند؟

در مدل‌های مبتنی بر ترنسفورمر، برخلاف مدل‌های RNN که به صورت ترتیبی داده‌ها را پردازش می‌کنند، تمام ورودی‌ها به طور همزمان پردازش می‌شوند. این موضوع باعث می‌شود که مدل اطلاعات Positional Embedding مربوط به ترتیب کلمات را به صورت ذاتی در اختیار نداشته باشد. به همین دلیل، یا کدگذاری موقعیتی در مدل‌های ترنسفورمر استفاده می‌شود تا اطلاعات مربوط به موقعیت هر کلمه در جمله به مدل اضافه شود و ساختار ترتیب کلمات حفظ گردد. این کار به مدل کمک می‌کند تا وابستگی‌های معنایی میان کلمات را به درستی تشخیص دهد و جملات منسجم‌تری تولید کند. روش Positional Embedding معمولاً به دو شکل انجام می‌شود: یکی آموزش‌پذیر (Learnable Positional Embedding) که در آن موقعیت‌های کلمات مانند خود واژگان به عنوان یک ماتریس وزنی یاد گرفته می‌شوند، و دیگری روش سینوسی و کسینوسی (Sinusoidal Positional Encoding) که در آن از توابع ریاضی برای تعیین مقادیر موقعیتی استفاده می‌شود. در روش سینوسی، هر موقعیت از یک سری مقادیر سینوسی و کسینوسی منحصر به فرد برای هر بعد استفاده می‌کند تا به مدل این امکان را بدهد که موقعیت نسبی کلمات را بدون نیاز به پارامترهای اضافی در نظر بگیرد. این نوع کدگذاری باعث می‌شود که مدل بتواند فاصله بین کلمات را درک کند و وابستگی‌های طولانی‌مدت را به خوبی حفظ کند. استفاده از Positional Embedding برای مدل ترنسفورمر ضروری است، زیرا بدون آن، مدل نمی‌تواند تفاوت بین جملاتی مانند "گربه روی میز است" و "میز روی گربه است" را متوجه شود. این روش اطلاعاتی درباره‌ی ترتیب و ساختار جمله را به مدل اضافه کرده و باعث می‌شود که مدل بتواند روابط بین کلمات را در نظر بگیرد و تولید متن دقیق‌تری داشته باشد. همچنین، این روش به مدل کمک می‌کند که وابستگی‌های معنایی بین کلمات را حتی در جملات طولانی حفظ کند، که یکی از نقاط ضعف مدل‌های RNN بود. در نهایت، Positional Embedding یک بخش حیاتی در مدل‌های ترنسفورمر است که بدون آن، مدل نمی‌تواند درک درستی از ترتیب کلمات داشته باشد و دقت پیش‌بینی کاهش پیدا می‌کند.

پس از اتمام آموزش، تغییرات Loss را برای داده‌های آموزش و اعتبارسنجی در قالب یک نمودار رسم کنید.



شکل ۱۷ نمودار لاس برای داده ترین و ولیدیشن

این نمودار روند کاهش خطای آموزش (Train Loss) و خطای اعتبارسنجی (Validation Loss) را در طول ۱۵ دوره‌ی آموزشی (Epoch) نمایش می‌دهد. در ابتدای آموزش، مقدار خطا در هر دو مجموعه بالا است که طبیعی بوده و نشان‌دهنده‌ی این است که مدل در مراحل اولیه هنوز یادگیری مناسبی نداشته و در حال تنظیم وزن‌های خود برای بهینه‌سازی پیش‌بینی‌ها است. مشاهده می‌شود که خطای آموزش به‌طور مداوم کاهش یافته که این امر بیانگر این است که مدل به درستی در حال یادگیری الگوهای موجود در داده‌ها است و به مرور تطبیق بهتری با مجموعه‌ی آموزشی پیدا کرده است. در ابتداء، خطای اعتبارسنجی نیز روندی مشابه دارد و همراه با خطای آموزش کاهش می‌یابد، اما پس از حدود ایپاک ۴، فاصله‌ی بین دو نمودار بیشتر شده و روند کاهش خطای اعتبارسنجی نسبت به خطای آموزش کندرتر می‌شود.

این تفاوت نشان‌دهنده‌ی آن است که مدل در حال یادگیری ویژگی‌های خاص مجموعه آموزشی است و احتمال بیش‌برازش (Overfitting) وجود دارد. در مراحل میانی آموزش، یعنی بین ایپاک ۶ تا ۱۰، مشاهده می‌شود که خطای آموزش همچنان کاهش می‌یابد اما خطای اعتبارسنجی تقریباً ثابت می‌ماند و حتی در برخی قسمت‌ها کمی نوسان دارد. این مسئله نشان می‌دهد که مدل ممکن است بیش از حد بر روی داده‌های آموزشی تمرکز کرده و توانایی تعمیم آن روی داده‌های جدید کمی کاهش یافته باشد. این وضعیت در ایپاک‌های انتهایی بیشتر آشکار می‌شود، جایی که خطای آموزش همچنان کاهش می‌یابد اما

خطای اعتبارسنجی تقریباً در یک مقدار مشخص باقی مانده و تغییر چندانی ندارد. این اتفاق معمولاً زمانی رخ می‌دهد که مدل به اندازه‌ی کافی آموزش دیده و ظرفیت بیشتری برای یادگیری ویژگی‌های عمومی‌تر ندارد، در حالی که هنوز برخی جزئیات خاص مجموعه‌ی آموزشی را یاد می‌گیرد که در داده‌های جدید دیده نخواهد شد.

برای کاهش اثر بیش‌برازش می‌توان از تکنیک‌هایی مانند Dropout، افزایش داده‌ها (Data Augmentation) و تنظیم وزن‌ها (Regularization) استفاده کرد. همچنین، یکی از راهکارهای مهم برای جلوگیری از بیش‌برازش در این مدل این است که آموزش را در نقطه‌ای که خطای اعتبارسنجی دیگر کاهش نمی‌یابد متوقف کنیم، که این کار می‌تواند از طریق Early Stopping انجام شود. در کل، این نمودار نشان می‌دهد که مدل به درستی آموزش دیده اما نیاز به تنظیم بیشتر برای بهبود تعمیم‌پذیری روی داده‌های جدید دارد تا دقت مدل در تولید کپشن‌ها به حداکثر ممکن برسد.

برای ۵ عکس تصادفی از مجموعه داده تست، کپشن تولید کنید و نتایج را بررسی کنید.

در این ۵ نمونه از داده‌های تست، مدل تلاش کرده است تا بر اساس اطلاعات تصویری ورودی کپشن‌هایی تولید کند، اما در برخی از موارد، اختلاف قابل توجهی میان کپشن تولیدشده (Pred) و کپشن واقعی (GT) مشاهده می‌شود که می‌تواند نشان‌دهنده‌ی نقاط ضعف مدل در درک برخی ویژگی‌های تصویر باشد. تحلیل موردنی این نمونه‌ها نشان می‌دهد که مدل در برخی موقعیت‌ها عملکرد خوبی داشته اما همچنان در بخش‌هایی دچار چالش است.

در نمونه اول، مدل "یک مرد و یک زن در حال راه رفتن روی پیاده‌رو هستند" را به عنوان خروجی تولید کرده، در حالی که کپشن واقعی تصویر نشان می‌دهد که "پسری از داخل یک چاله‌ی گل‌آلود می‌پردد". این اختلاف نشان می‌دهد که مدل در تشخیص دقیق فعالیت‌های سوژه‌ها مشکل دارد و ممکن است به دلیل یادگیری وابستگی‌های قوی از مجموعه داده‌ی آموزشی، به جای تمرکز روی حرکات و رفتارهای خاص، از الگوهای رایج‌تر در داده‌های آموزشی استفاده کند.

در نمونه دوم، مدل کپشن "یک مرد و یک زن روی یک نیمکت نشسته‌اند" را تولید کرده در حالی که کپشن واقعی نشان می‌دهد "یک زوج به همراه نوزادشان زیر یک درخت روی دریاچه نشسته‌اند". در این مورد، مدل مفهوم کلی نشستن افراد را درست تشخیص داده اما حضور نوزاد و دریاچه را نادیده گرفته است که نشان می‌دهد مدل در تشخیص اشیای کوچک‌تر و پس‌زمینه‌ها دقت کافی ندارد.

در نمونه سوم، مدل تشخیص داده که "مردی با یک پیراهن مشکی در حال انجام یک حرکت نمایشی روی اسکیت‌برد است"، در حالی که کپشن واقعی تصویر به "یک دوچرخه‌سوار BMX که در حال انجام

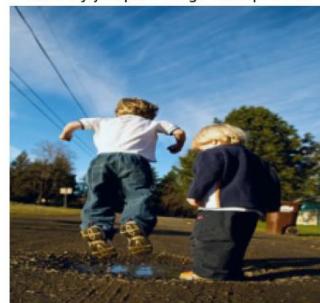
حرکت معلق در هوا است" اشاره دارد. این خطا نشان می‌دهد که مدل در تشخیص تفاوت بین دوچرخه و اسکیت‌برد دچار اشتباه شده که ممکن است به دلیل شباهت ساختاری این دو وسیله در برخی زوایای تصویر باشد.

در نمونه چهارم، مدل کپشن "یک اسکیت‌باز در حال انجام یک حرکت نمایشی روی یک رمپ است" را تولید کرده که تا حدی به کپشن واقعی که "یک اسکیت‌باز در هوا مقابل ساختمان‌های کوتاه آجری است" نزدیک است اما تفاوت‌هایی نیز وجود دارد. در این مورد، مدل مفهوم کلی حرکت نمایشی اسکیت‌باز را به درستی تشخیص داده اما جزئیات محیطی و پس‌زمینه را در نظر نگرفته که می‌تواند ناشی از ضعف مدل در یادگیری ارتباطات پس‌زمینه‌ای باشد.

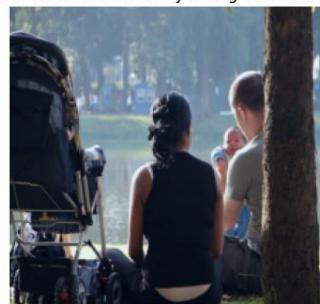
در نمونه پنجم، مدل کپشن "یک سگ در حال بازی با یک توپ روی چمن است" را تولید کرده که با کپشن واقعی "یک سگ بزرگ بالای سر یک سگ کوچک‌تر ایستاده که دهانش باز است" تفاوت دارد. این خطا نشان می‌دهد که مدل در تشخیص تعامل بین چندین شیء (در اینجا دو سگ) عملکرد مناسبی نداشت و به اشتباه فقط بر روی توپ تمرکز کرده است که می‌تواند ناشی از تمایل مدل به استفاده از اطلاعات پرتکرار در مجموعه آموزشی باشد.

به طور کلی، بررسی این نمونه‌ها نشان می‌دهد که مدل توانایی تشخیص سوژه‌های اصلی را دارد اما در تفکیک جزئیات، فعالیت‌های خاص و روابط میان اشیا دچار چالش است. این مشکلات را می‌توان با استفاده از داده‌های آموزشی متنوع‌تر، بهبود مکانیزم توجه (Attention Mechanism) و استفاده از روش‌های پیشرفته‌تر مانند ترنسفورمرهای بینایی (Vision Transformers) و مدل‌های چندوجهی (Multimodal Models) کاهش داد تا دقیق‌تر افزایش یابد.

Pred: a man and a woman are walking on a sidewalk
GT: boy jumps through mud puddle



Pred: a man and a woman are sitting on a bench
GT: a couple with their newborn baby sitting under a tree facing a lake



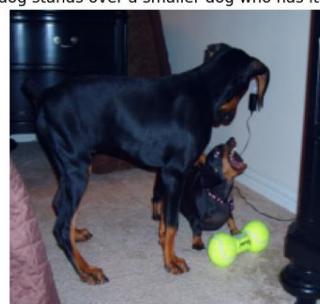
Pred: a man in a black shirt is doing a stunt on a skateboard
GT: a bmx biker is currently in midair doing a front flip



Pred: a skateboarder does a trick on a ramp
GT: a skateboarder is airborne in front of some low brick buildings



Pred: a dog is playing with a ball in the grass
GT: a large dog stands over a smaller dog who has its mouth open



شکل ۱۸ نمونه کپشن های تولید شده

[INFO] Showing 5 sample predictions from test set.

[INFO] Evaluation on Validation dataset:

BLEU-1: 0.3760 | BLEU-2: 0.2157 | BLEU-3: 0.1253 | BLEU-4: 0.0696
ROUGE-L: 0.3741

[INFO] Evaluation on Test dataset:

BLEU-1: 0.3708 | BLEU-2: 0.2139 | BLEU-3: 0.1300 | BLEU-4: 0.0753
ROUGE-L: 0.3733

[INFO] All steps completed.

[TIME] 372.62 seconds total run time.

نتایج ارزیابی مدل روی مجموعه داده‌های اعتبارسنجی و تست نشان می‌دهند که مدل عملکرد قابل قبولی در تولید کپشن‌های مرتبط با تصاویر دارد اما همچنان محدودیت‌هایی در دقت پیش‌بینی وجود دارد. شاخص BLEU که معیار اندازه‌گیری شباهت بین جملات تولیدی و جملات واقعی است، نشان می‌دهد که مدل در تشخیص کلمات منفرد و توالی‌های کوتاه عملکرد بهتری دارد اما با افزایش طول توالی‌ها، دقت مدل کاهش پیدا می‌کند. مقدار BLEU-1 برابر ۰,۳۷۶۰ در مجموعه اعتبارسنجی و ۰,۳۷۰۸ در مجموعه تست است که نشان می‌دهد مدل در سطح کلمات منفرد و ارتباط آن‌ها با کپشن واقعی، دقت مناسبی دارد اما این مقدار با کاهش وزن کلمات منفرد و افزایش تمرکز روی n-gram‌های طولانی‌تر، افت پیدا می‌کند. مقدار BLEU-2 که نشان‌دهنده دقت مدل در سطح دوکلمه‌ای است، برابر ۰,۲۱۵۷ در مجموعه اعتبارسنجی و ۰,۲۱۳۹ در مجموعه تست است که کاهش نسبتاً قابل توجهی نسبت به BLEU-1 دارد و بیانگر این است که مدل در درک ارتباط بین کلمات متوالی هنوز نیاز به بهبود دارد. مقدار BLEU-3 برابر ۰,۱۲۵۳ در اعتبارسنجی و ۰,۱۳۰۰ در تست نشان می‌دهد که مدل در یادگیری الگوهای طولانی‌تر همچنان دچار مشکل است. مقدار BLEU-4 که سخت‌گیرانه‌ترین معیار است، برابر ۰,۶۹۶۰ در اعتبارسنجی و ۰,۰۷۵۳ در تست است که نشان‌دهنده این است که مدل در تولید جملات کاملاً مشابه با کپشن‌های واقعی همچنان عملکرد ضعیفی دارد و اغلب کپشن‌های تولیدی از لحاظ ساختاری و ترتیب کلمات تفاوت‌هایی با کپشن‌های مرجع دارند.

شاخص ROUGE-L که بر اساس میزان همپوشانی طولانی‌ترین زیرشته‌های مشترک بین جمله واقعی و جمله تولیدی محاسبه می‌شود، مقدار ۰,۳۷۴۱ در اعتبارسنجی و ۰,۳۷۳۳ در تست را نشان می‌دهد که بیانگر آن است که مدل در تولید جملات با ساختار مشابه کپشن واقعی تا حدی موفق بوده اما همچنان برخی از کپشن‌های تولیدی از لحاظ ترتیب و چینش کلمات با کپشن‌های اصلی تفاوت دارند.

در مجموع، این نتایج نشان می‌دهند که مدل توانایی درک و توصیف کلی تصاویر را دارد اما در حفظ ترتیب دقیق کلمات و تولید جملات طولانی‌تر و روان‌تر دچار چالش است. برای بهبود عملکرد مدل می‌توان از Sampling به جای Beam Search ساده، استفاده از داده‌های آموزشی متنوع‌تر، و اعمال مکانیزم‌های

بهبود یافته‌ی توجه (Self-Attention Mechanisms) بهره برد. همچنین، استفاده از مدل‌های پیشرفته‌تر مانند ترنسفورمرهای چندوجهی و مدل‌های Vision-Language مانند Flamingo و BLIP می‌تواند دقت مدل را در تولید کپشن‌های غنی‌تر و معنادارتر افزایش دهد.

۶. امتیازی

تحقیق کنید چه معیارهایی برای ارزیابی این مدلها بکار میرود

ارزیابی مدل‌های Image Captioning نیازمند معیارهایی است که بتوانند کیفیت کپشن‌های تولیدشده را نسبت به کپشن‌های واقعی سنجیده و نشان دهنند که مدل تا چه حد توانسته است توصیفی صحیح، روان و شبیه به زبان انسانی ارائه دهد. این معیارها به دو دسته اصلی تقسیم می‌شوند: معیارهای مبتنی بر شbahت متنی و معیارهای مبتنی بر ارزیابی انسانی. معیارهای مبتنی بر شbahت متنی شامل BLEU، SPICE، METEOR، ROUGE و CIDEr هستند. BLEU میزان همپوشانی n-gram‌ها را بین کپشن تولیدی و کپشن واقعی بررسی می‌کند و هرچه n-gram طولانی‌تر باشد، معیار سخت‌گیرانه‌تر می‌شود. مقدار BLEU-1 شباهت تک‌کلمه‌ای، BLEU-2 دوکلمه‌ای، BLEU-3 سه‌کلمه‌ای و BLEU-4 چهار‌کلمه‌ای را بررسی می‌کند. این معیار اگرچه پرکاربرد است اما فقط بر اساس تطابق سطحی کلمات کار می‌کند و معنا و مفهوم کلی را در نظر نمی‌گیرد. ROUGE نیز میزان شباهت بین کپشن‌های تولیدی و کپشن‌های واقعی را بر اساس طولانی‌ترین زیرشته‌های مشترک (LCS) محاسبه می‌کند و بیشتر روی Recall تمرکز دارد. METEOR که مکمل BLEU است، علاوه بر همپوشانی واژگان، ریشه‌یابی کلمات، متراffد‌ها و ترتیب جملات را در نظر می‌گیرد و برای ارزیابی شباهت متنی دقت بیشتری دارد. CIDEr یکی دیگر از معیارهای مهم ارزیابی است که با استفاده از روش TF-IDF، وزن بیشتری به کلمات مهم و SPICE نادر اختصاص می‌دهد و کپشن‌هایی که شامل واژگان کلیدی باشند امتیاز بالاتری دریافت می‌کنند. یک معیار پیشرفته‌تر است که به جای تمرکز صرف بر شباهت واژگانی، روابط معنایی بین کلمات و ساختار جمله را تحلیل می‌کند و از گراف معنایی برای ارزیابی میزان درک مدل از اجزای تصویر استفاده می‌کند. در کنار معیارهای خودکار، معیارهای مبتنی بر ارزیابی انسانی نیز بسیار مهم هستند. Fluency یا روان بودن جمله بررسی می‌کند که آیا کپشن تولید شده از لحاظ دستوری و زبانی صحیح است یا نه. Relevance میزان ارتباط کپشن با محتوای تصویر را ارزیابی می‌کند. Diversity مشخص می‌کند که آیا مدل قادر به تولید جملات متنوع و غیرتکراری است یا خیر. Consistency بررسی می‌کند که آیا توصیف تولیدی سازگاری منطقی با تمامی اجزای تصویر دارد. هر کدام از این معیارها نقاط ضعف و قوت خود را دارند. SPICE، ROUGE و METEOR بیشتر به همپوشانی زبانی توجه دارند، درحالی که CIDEr و BLEU علاوه بر متن، به مفاهیم معنایی و اهمیت کلمات نیز دقت می‌کنند. SPICE برای بررسی معنا و ارتباط

منطقی اجزای جمله با تصویر عملکرد بهتری دارد، در حالی که CIDEr برای تشخیص کلمات کلیدی بهتر عمل می‌کند. ارزیابی انسانی نیز از آنجا که درک واقعی از عملکرد مدل را ارائه می‌دهد، مکمل این معیارها محسوب می‌شود. در مجموع، بهترین روش برای ارزیابی یک مدل Image Captioning ترکیب چند معیار مانند BLEU، SPICE، CIDEr و ارزیابی انسانی است تا علاوه بر شباهت متنی، مفهوم و کیفیت واقعی کپشن‌های تولیدی نیز سنجیده شود.

Bleu Score انجام داده و بهطور مختصر توضیح دهدید که این معیار چگونه عملکرد مدل را ارزیابی می‌کند. سپس (Bleu Score از یک تا چهار) را روی داده‌های تست محاسبه کرده و نتایج آن را برای بخش‌های مختلف گزارش کنید و مقایسه‌های انجام دهید.

از جمله مدل‌های Image Captioning است. این معیار به طور خاص برای سنجش کیفیت ترجمه‌های ماشینی توسعه یافته است اما در بسیاری از مسائل مرتبط با پردازش زبان طبیعی (NLP) مانند خلاصه‌سازی متن و تولید توضیحات تصویری نیز کاربرد دارد. BLEU میزان شباهت بین یک جمله‌ی تولید شده توسط مدل (Hypothesis) و یک یا چند جمله‌ی مرجع (Reference) را اندازه‌گیری می‌کند و این شباهت را بر اساس میزان همپوشانی n-gram-ها بین متن‌های تولیدی و متن‌های مرجع محاسبه می‌کند.

روش کار BLEU Score به این صورت است که ابتدا n-gram-های مشترک بین خروجی مدل و جملات مرجع را شمارش کرده و یک نسبت (Precision) از تعداد n-gram-های منطبق را به دست می‌آورد. مقدار BLEU-2 می‌تواند از ۱ تا ۴ باشد که به ترتیب شامل (Unigram)، (BLEU-1 تطابق تک کلمه‌ای یا (Trigram)، (BLEU-3 تطابق سه کلمه‌ای یا (Bigram) و (BLEU-4 تطابق دو کلمه‌ای یا (چهار کلمه‌ای) است. هرچه مقدار n بزرگ‌تر باشد، BLEU میزان سخت‌گیرانه‌تری برای ارزیابی کیفیت متن تولیدی خواهد بود، زیرا علاوه بر دقیقت کلمات منفرد، ترتیب و ساختار جملات را نیز بررسی می‌کند.

مراحل محاسبه BLEU Score:

محاسبه n-gram Precision: ابتدا تعداد n-gram-های مشترک بین کپشن تولید شده و کپشن‌های مرجع محاسبه شده و نسبت آن به کل n-gram-های تولید شده مشخص می‌شود. به عنوان مثال، اگر مدل جمله‌ی "A brown dog is playing" را تولید کرده باشد و کپشن مرجع "A brown dog is running" باشد، که ۵ مورد از ۵ n-gram-های مشترک برای BLEU-1 عبارتند از {"A", "brown", "dog", "is"} کلمه را شامل می‌شود، بنابراین Precision برابر ۵/۴ خواهد بود.

اعمال پنالتی برای طول جملات (Precision - BP): که مدل ممکن است با تولید جملات کوتاه، دقت بالایی کسب کند. برای جلوگیری از این مشکل، یک جریمه‌ی طول اضافه می‌کند تا جملات بسیار کوتاه، امتیاز غیرواقعی بالایی دریافت نکنند. اگر کپشن تولیدی از کپشن مرجع کوتاه‌تر باشد، BLEU از یک ضریب جریمه استفاده می‌کند که مقدار آن به صورت نمایی از نسبت طول جمله تولید شده به جمله مرجع محاسبه می‌شود.

میانگین هندسی BLEU (Geometric Mean): برای ترکیب n-gram‌های مختلف، به جای میانگین حسابی، از میانگین هندسی Precision‌های مختلف استفاده می‌کند که باعث می‌شود مقادیر کوچکتر تأثیر بیشتری داشته باشند. به عبارت دیگر، اگر BLEU-1 مقدار بالایی داشته باشد اما BLEU-4 مقدار پایینی داشته باشد، میانگین هندسی تأثیر BLEU-4 را افزایش داده و مدل را به تولید جملات کامل‌تر و طبیعی‌تر تشویق می‌کند. فرمول آن به نحو زیر است:

$$BLEU = BP \times \exp \left(\sum_{n=1}^N w_n \log P_n \right)$$

مزایا و محدودیت‌های BLEU Score: این است که یک معیار سریع، محاسباتی و مستقل از زبان است و نیازی به ارزیابی انسانی ندارد. این معیار به‌طور گسترده در پژوهش‌های مختلف استفاده می‌شود و امکان مقایسه مدل‌های مختلف را فراهم می‌کند. با این حال، محدودیت‌هایی نیز دارد. اول اینکه BLEU فقط به همپوشانی کلمات توجه می‌کند و ساختار معنایی جمله و روابط بین اجزا را در نظر نمی‌گیرد. به همین دلیل ممکن است یک جمله که از لحاظ معنایی درست است اما با کلمات متفاوتی بیان شده است، امتیاز BLEU پایینی بگیرد. همچنین، این معیار برای جملات کوتاه و خلاقانه مانند کپشن تصاویر چندان دقیق نیست، زیرا ممکن است چندین جمله معتبر برای توصیف یک تصویر وجود داشته باشد که در کپشن‌های مرجع ذکر نشده است.

این معیار برای همه مدل‌ها محاسبه شده که در بالا گزارش شده است. مقایسه اعداد بدست آمده از هر مدل در زیر آمده است:

برای مقایسه عملکرد مدل‌های مختلف در تولید کپشن تصاویر، سه مجموعه از BLEU Score‌ها را داریم که میزان دقت مدل‌ها در شباهت خروجی تولیدی با کپشن‌های واقعی را بررسی می‌کنند. ابتدا به تحلیل عملکرد هر مدل پرداخته و سپس نتیجه‌گیری خواهیم کرد که کدام مدل عملکرد بهتری داشته است. مدل LSTM که بر اساس CNN-RNN پیاده‌سازی شده، نتایج ضعیفتری را نشان می‌دهد. مقدار BLEU-1

برابر با ۰،۳۰۲۰ و BLEU-4 مقدار بسیار پایین ۰،۰۳۶۸ دارد که نشان می‌دهد مدل در تشخیص تکوازه‌ها عملکرد متوسطی دارد اما در درک ترتیب کلمات و تولید جملات بلند دچار مشکل است. دلیل این ضعف به ماهیت LSTM برمی‌گردد که وابستگی‌های طولانی‌مدت را به خوبی یاد نمی‌گیرد و برای تولید جملات بلندتر دچار کاهش دقت می‌شود. مدل CNN-RNN همراه با مکانیزم Attention توانسته است عملکرد بهتری نسبت به مدل قبلی داشته باشد. مقدار BLEU-1 در این مدل به ۰،۳۶۳۵ و BLEU-4 به ۰،۰۸۳۲ افزایش یافته است که نشان‌دهنده بهبود در درک وابستگی‌های بین کلمات و همچنین هماهنگی بهتر جملات تولیدی است. مکانیزم توجه باعث شده مدل بتواند روی بخش‌های مهم تصویر تمرکز کند و هنگام تولید هر واژه، ارتباط آن را با قسمت‌های مختلف تصویر برقرار کند. با این حال، همچنان BLEU-4 پایین است که نشان می‌دهد مدل در تولید توصیف‌های طولانی‌تر هنوز محدودیت دارد. در نهایت، مدل-CNN-Transformer با مقدار BLEU-1 برابر با ۰،۳۷۶ و BLEU-2 برابر با ۰،۲۱۵۷ عملکرد بهتری در انتخاب کلمات مناسب داشته است اما مقدار BLEU-3 و BLEU-4 نسبت به مدل Attention-based کمی پایین‌تر است که نشان می‌دهد مدل در ساختار جمله‌بندی نیاز به بهبود دارد. دلیل این موضوع این است که ترنسفورمرها در پردازش ترتیب جملات نسبت به مکانیزم Attention بهینه‌تر هستند اما اگر داده‌های آموزشی و تعداد epoch‌ها افزایش یابند، می‌توانند بهبود چشمگیری داشته باشند. در مقایسه کلی، مدل CNN-Transformer در انتخاب واژگان دقت بهتری دارد، اما مدل CNN-RNN همراه با مکانیزم Attention در حفظ ترتیب کلمات و درک ساختار جملات طولانی‌تر عملکرد بهتری نشان داده است. مدل LSTM ساده به دلیل نبود مکانیزم توجه و ضعف در نگهداری اطلاعات بلندمدت، کمترین دقت را در بین مدل‌های مقایسه شده دارد. برای بهبود این مدل‌ها، می‌توان از روش‌هایی مانند افزایش حجم داده‌های آموزشی، استفاده از مکانیزم Self-Attention (ViTs) پیشرفته‌تر مانند Vision Transformers (ViTs) و همچنین بهره‌گیری از مدل‌های چندوجهی (Multimodal Flamingo) یا BLIP استفاده کرد. در نهایت، بهترین مدل بسته به کاربرد موردنظر انتخاب می‌شود؛ اگر هدف دقت بالاتر در انتخاب کلمات باشد، مدل Transformer مناسب‌تر است اما اگر هدف بهبود در ترتیب جملات و درک بهتر ساختار کپشن‌ها باشد، مدل Attention-based RNN گزینه بهتری خواهد بود.